

Petascale Computing in the U.S.

Horst D. Simon

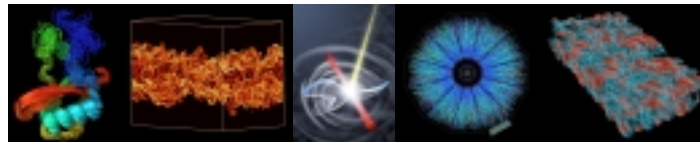
Associate Laboratory Director, Computing Sciences
Lawrence Berkeley National Laboratory

**ACTS Workshop
Berkeley, California
August 24, 2006**



Acknowledgements

- **Many colleagues contributed slides, ideas, and valuable feedback:**
 - **Earl Joseph (IDC), Jeremy Kepner (MIT), Bill Kramer (NERSC), Steve Meacham (NSF), Randy Moulic (IBM), Dave Patterson (UC Berkeley), John Shalf (NERSC), David Skinner (NERSC), Rick Stevens (ANL), Erich Strohmaier (LBNL), Pete Ungaro (Cray), Howard Walter (NERSC)**
 - **Participants of Dagstuhl Workshop on Petascale Algorithms and Applications, 2006**
 - **TOP500 team: Erich Strohmaier, Hans Meuer, Jack Dongarra**





National Energy Research Scientific Computing Center

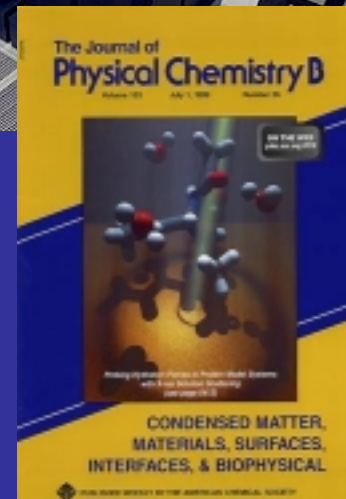
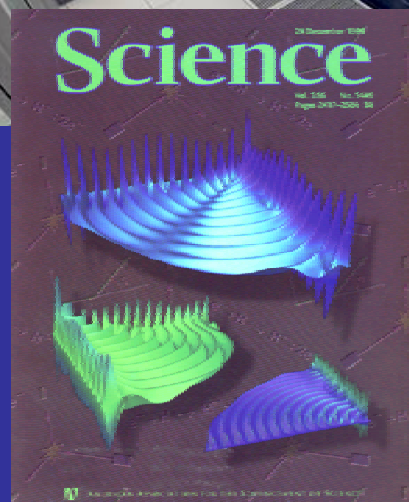
Serves the entire scientific
community

~2500 Users in
~250 projects

- Focus on
large-scale
computing



ERSC



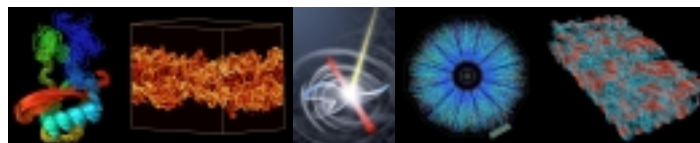
NERSC Center Overview

- **Funded by DOE, annual budget \$38M (FY06), about 60 staff**
 - Traditional strategy to invest equally in newest compute platform, staff, and other resources
- **Supports open, unclassified, basic research**
- **Close collaborations between university and NERSC in computer science and computational science**



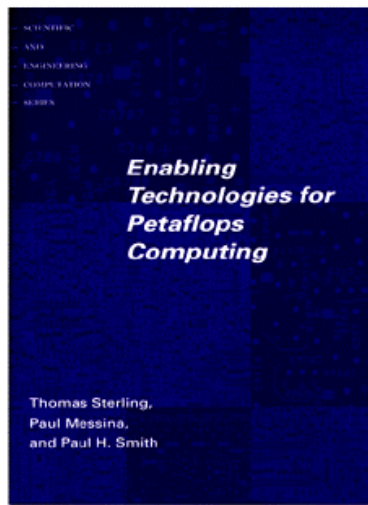
Overview

- **History and Future of Petaflops Computing**
- **HPC in 2006: “It was the best of times, it was the worst of times ...”**
- **“A Petaflops before its Time”**
- **The power problem**
- **The scaling problem**
- **What’s next?**

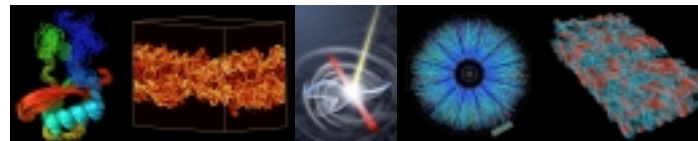


Steps Leading to Petaflops Computing

- Several workshops starting in the mid 1990s
 - 1994 Petaflops I (Pasadena)
 - 1995/1996 Summer Study (Bodega)
- 2002 DARPA HPCS
- 2003 HECRTF Roadmap
- 2004 NAS Report “The Future of Supercomputing”
- 2006 ACI (American Competitiveness Initiative)



1994



2006





DARPA HPCS Challenges

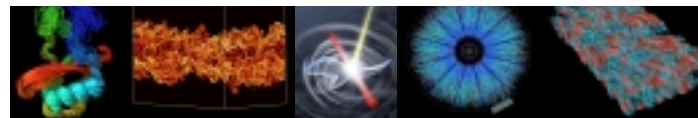
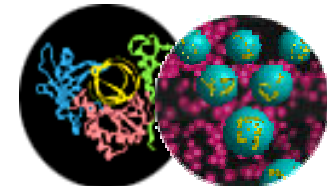
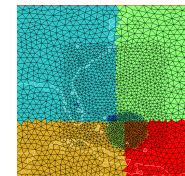
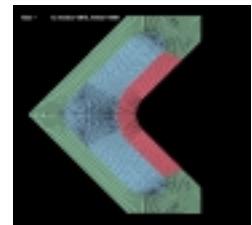
HPCS 

Goal:

- Provide a new generation of economically viable high productivity computing systems for the national security and industrial user community (2010)

Focus on:

- Real (not peak) performance of critical national security applications
 - Intelligence/surveillance
 - Reconnaissance
 - Cryptanalysis
 - Weapons analysis
 - Airborne contaminant modeling
 - Biotechnology
- Programmability: reduce cost and time of developing applications
- Software portability and system robustness





HPCS Roadmap

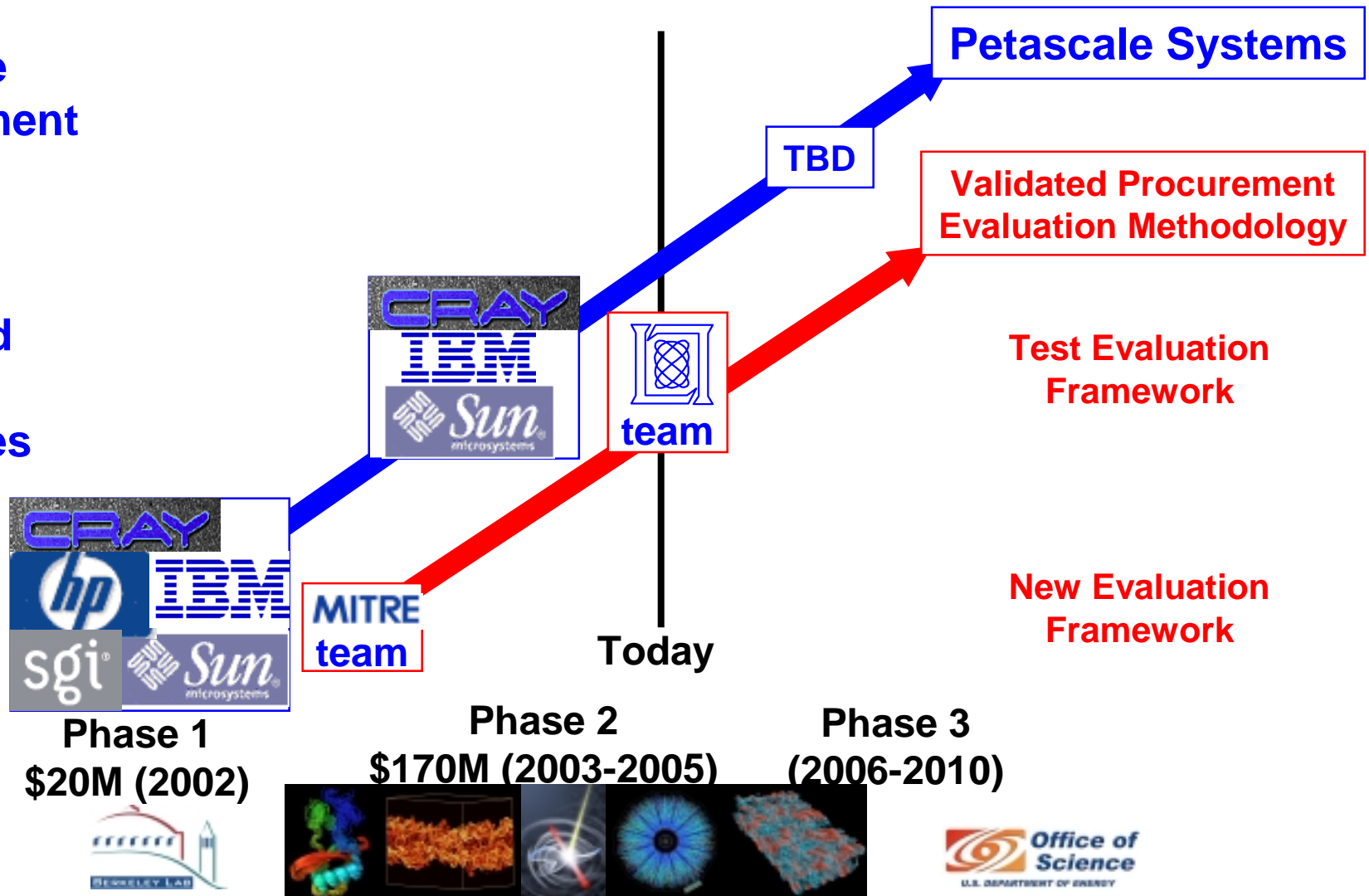
HPCS

- 5 vendors in phase 1; 3 vendors in phase 2; 1+ vendors in phase 3
- MIT Lincoln Laboratory leading measurement and evaluation team

Full Scale
Development

Advanced
Design &
Prototypes

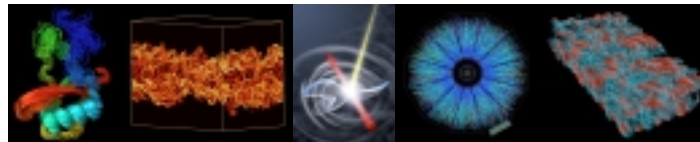
Concept
Study





DOE Office of Science Leadership Computing Facility Strategy

- **DOE selected ORNL and ANL to develop the DOE SC Leadership Computing Facilities**
 - ORNL will develop a series of systems based on Cray's XT3 and XT4 architectures with systems @ 250TF/s in FY07 and @1000TF/s in FY08/FY09
 - ANL will develop a series of systems based on IBM's BlueGene @ 100TF/s in FY07 and up to 1000TF/s in FY08/FY09 with BG/P
 - The Leadership Class Computing (LCC) systems are likely to be the most powerful civilian systems in the world when deployed
- **DOE SC will make these systems available as capability platforms to the broad national community via competitive awards (e.g. INCITE and LCC Allocations)**
 - Each facility will target ~20 large-scale production applications teams
 - Each facility will also support order 100 development users
- **DOE's LCC facilities will complement the existing and planned production resources at NERSC**
 - NERSC continues to support the broad based computational needs of the DOE SC mission
 - NERSC-5 will be at 100 TF/s in FY07 and NERSC-6 at 500 TF/s in FY10

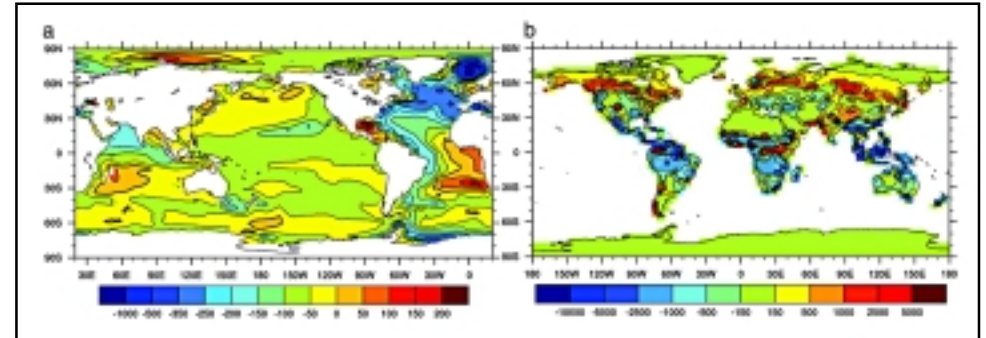


SciDAC - first federal program to implement Computational Science and Engineering

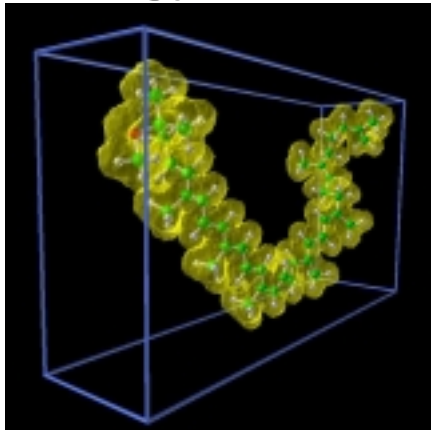
SciDAC (Scientific Discovery through Advanced Computing)

- About \$50M annual funding (2001 - 2006)
- 5 year continuation starting in FY07
- Focus on petascale applications FY07 - FY11

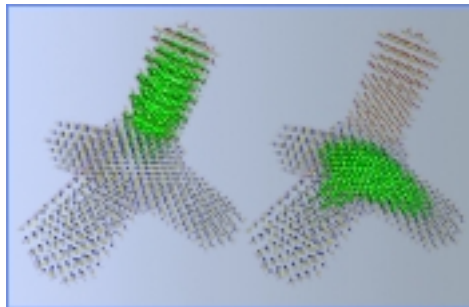
Global Climate



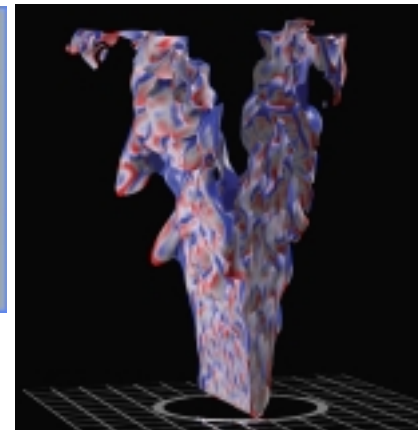
Biology



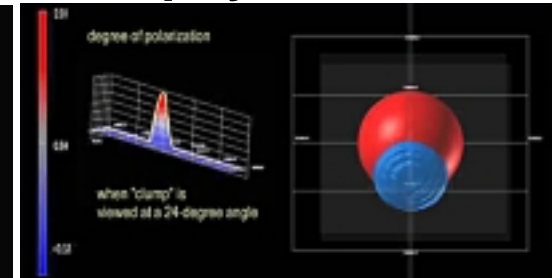
Nanoscience



Combustion

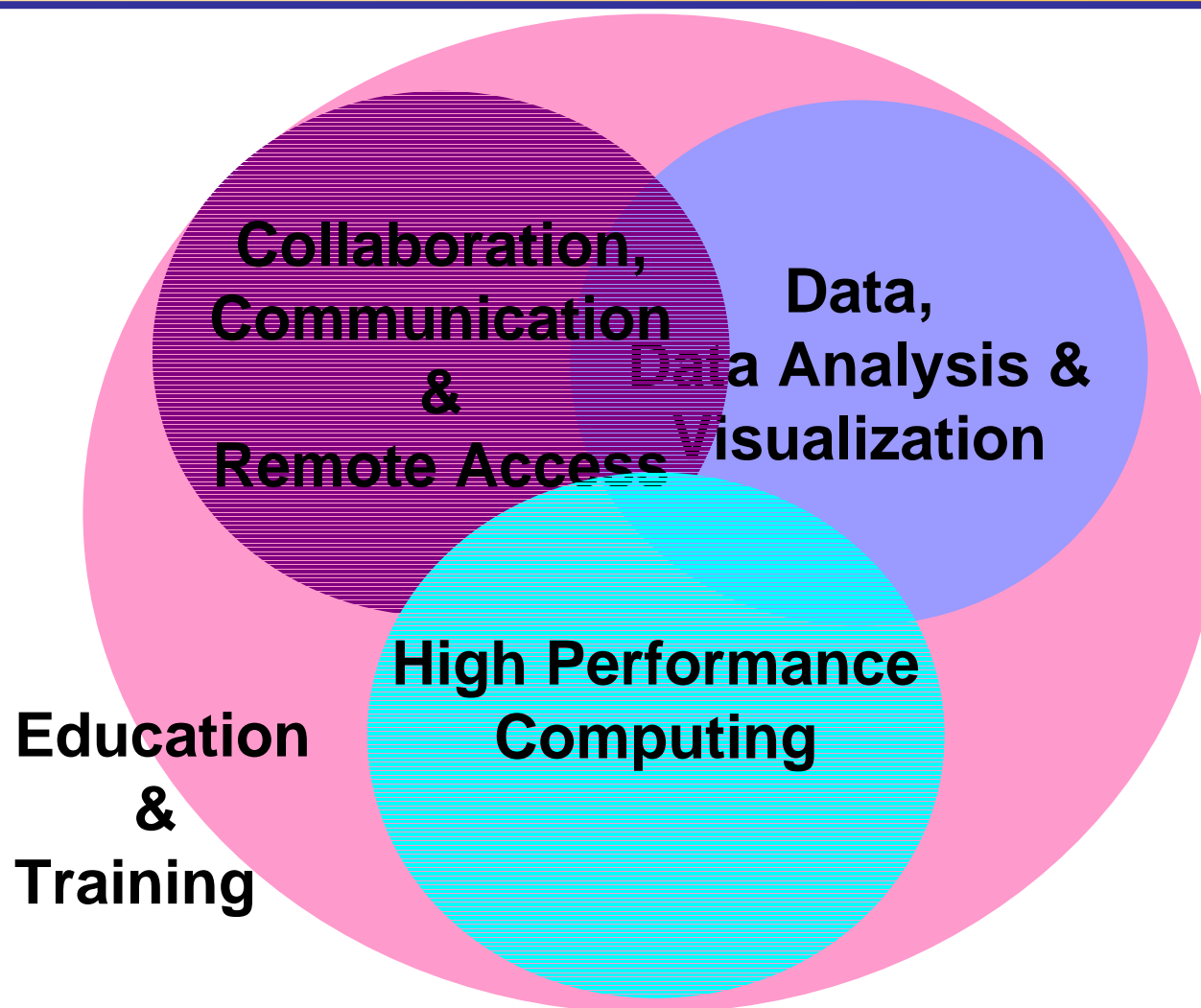


Astrophysics

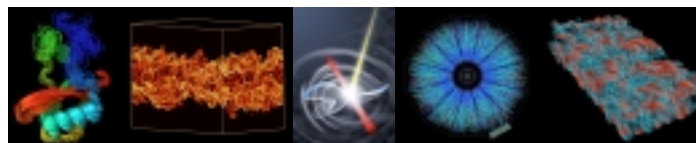




Office of Cyberinfrastructure (CI)

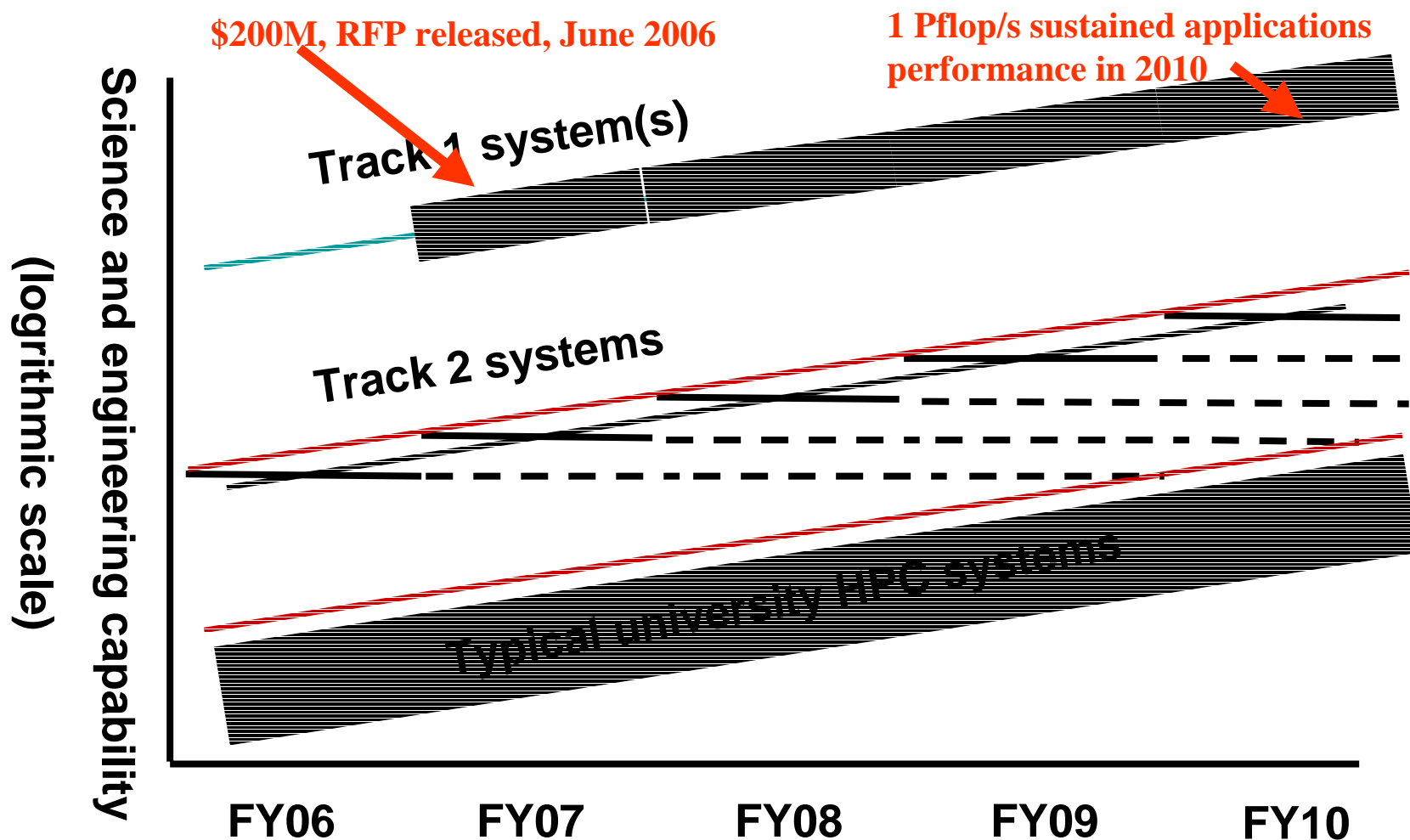


From Steve Meacham, NSF

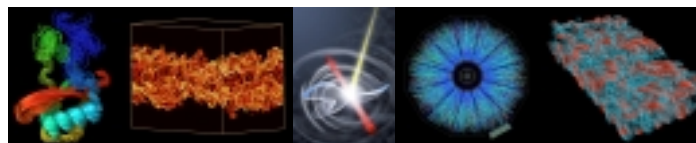




NSF Acquisition Strategy

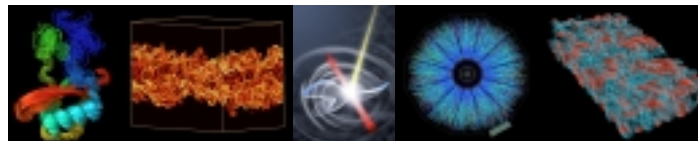


From Steve Meacham, NSF

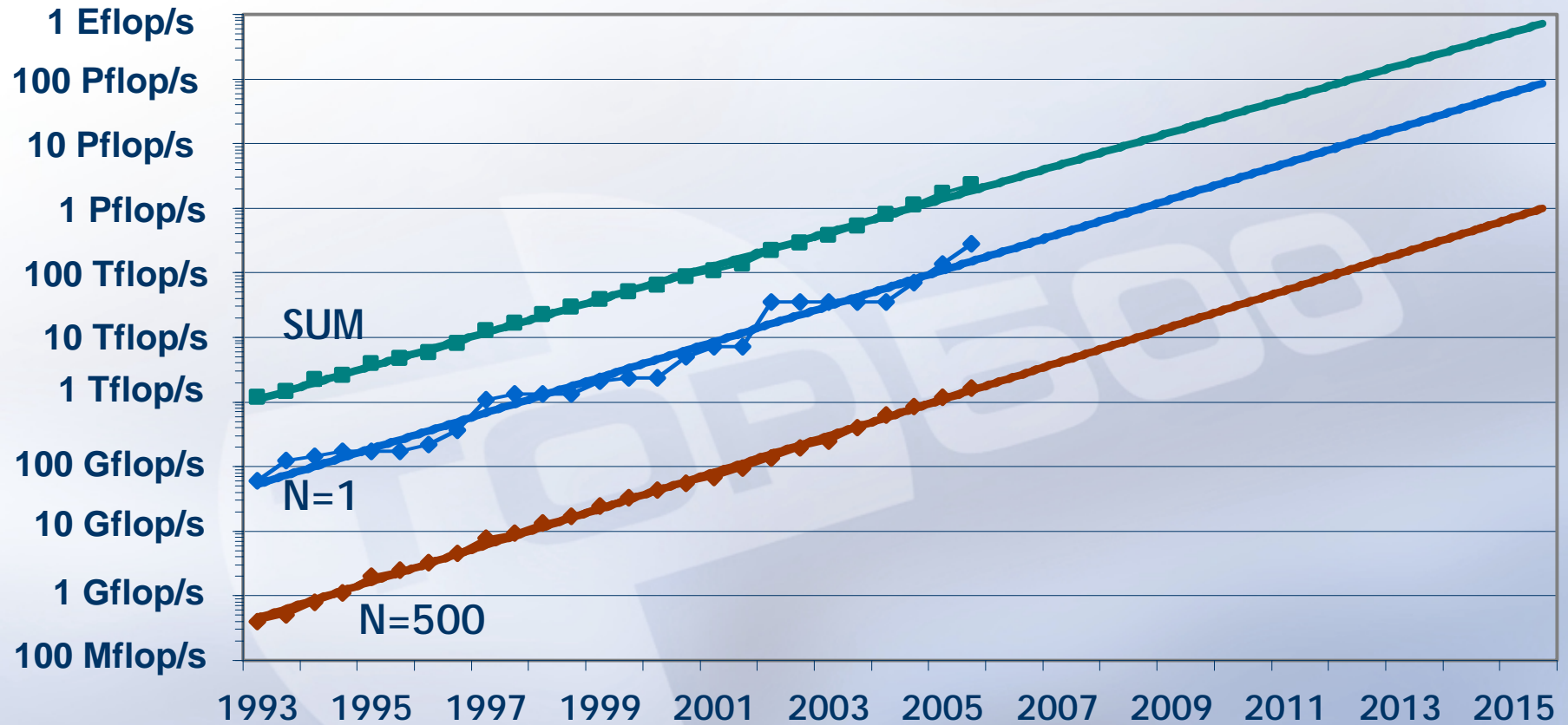


Levels of Petascale Computing

- The term “Petascale” is frequently used, but unfortunately ill-defined
- We need to distinguish
 - Theoretical peak petaflop/s systems
 - LINPACK Rmax Petaflop/s systems (used in TOP500)
 - Sustained applications performance in excess of a Petaflop/s
- My Definition: “**Petascale Computing**”
 - Widespread use of systems that deliver sustained applications performance a level above 1 Petaflop/s
 - Reached when all system on the TOP500 list have more than 1 Petaflop/s Rmax performance

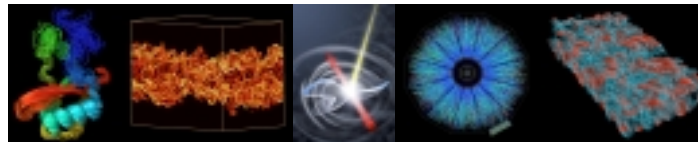


Performance Projection



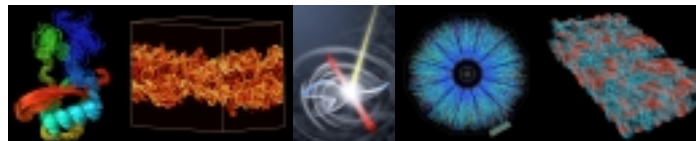
TOP500 Projections

- June 1997:
 - First LINPACK Teraflop/s system tops the list
- June 2005 (8 years later): Terascale computing arrives
 - 1Teraflop/s is required to enter the TOP500 list
- November 2008:
 - First LINPACK Petaflop/s system tops the list
- June 2016 (7.5 years later): Petascale computing arrives
 - 1 Petaflop/s is required to enter the TOP500 list
- November 2018:
 - First LINPACK Exaflop/s



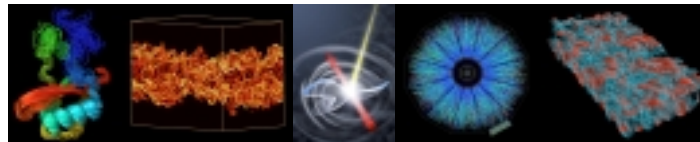
Platforms Capable of peak Petaflop/s in 2009

	Scale Demonstrated Factor to PF	Failures per Month Per TF	Power Consumption @ PF	Estimated System Cost
Cray XT3/XT4	10880 cpus 10x to PF ~100,000 cpus	~.1 - ~1	~8MW ^{XT4}	>\$150M ^{XT4} +memory
IBM Power5/6	10240 cpus 7x to PF ~72,000 cpus	1.3	~9.4MW ^{P6}	>\$170M ^{P6} +memory
Clusters x86-64/AMD64	8000 cpus 12x to PF ~100,000 cpus	2.6-8.0	~6MW ^{x86QC}	> \$150M ^{x86} +memory
Blue Gene L/P	131,720 cpus 2.2x to PF 294,912 cpus	.01-.03	~2.3MW ^P	< \$100M ^{BG} Including 288TB



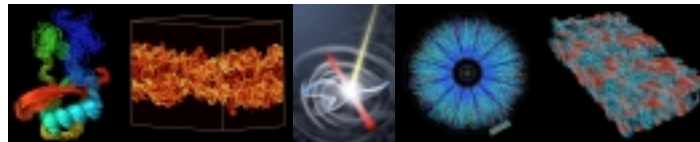
Overview

- **History and Future of Petaflops Computing**
- **HPC in 2006: “It was the best of times, it was the worst of times ...”**
- **“A Petaflops before its Time”**
- **The power problem**
- **The scaling problem**
- **What’s next?**



The Best of Times in 2005/2006

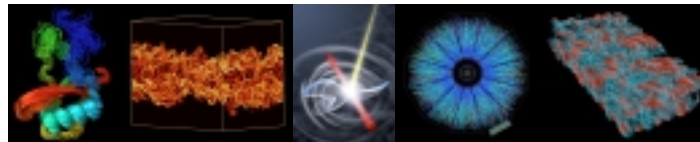
- HPC market continues to grow with double digit growth rates
- SC2005 record attendance indicator of optimism in HPC
- Increased political support for supercomputing
 - American Competitiveness Initiative (ACI)
 - Budget Increases for NSF, DOE/SC
 - “Petascale” everywhere (Japan “keisoku”, European Petascale Center)



The Worst of Times in 2005/2006

Customer perspective:

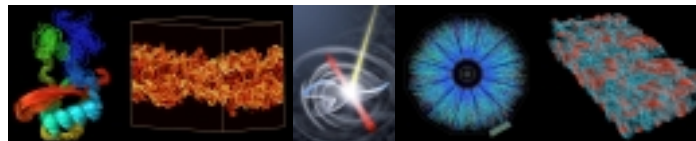
- **Traditional HPC vendors are struggling or are fading from market**
 - SGI in Chapter 11
 - Cray moving towards profitability
 - HP and Sun are not a strong presence at the high end
- **Newcomers are not ready yet for the high end**



IDC's HPC Market Definitions

- **Technical Capability**
 - Systems configured and purchased to solve the largest most demanding problems
- **Technical Enterprise**
 - Systems purchased to support technical applications in throughput environments selling for \$1 million or more
- **Technical Divisional**
 - Systems purchased for throughput environments selling from \$250,000 to \$999,000
- **Technical Departmental**
 - Systems purchased for throughput environments selling for \$50,000 to \$250,000
- **Technical Workgroup**
 - Systems under \$50,000

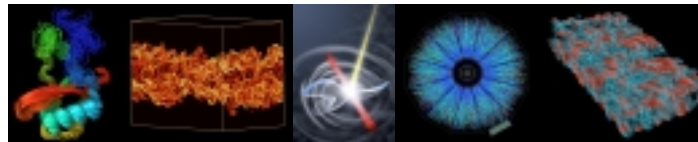
Slide courtesy of Earl Joseph Jr. , IDC



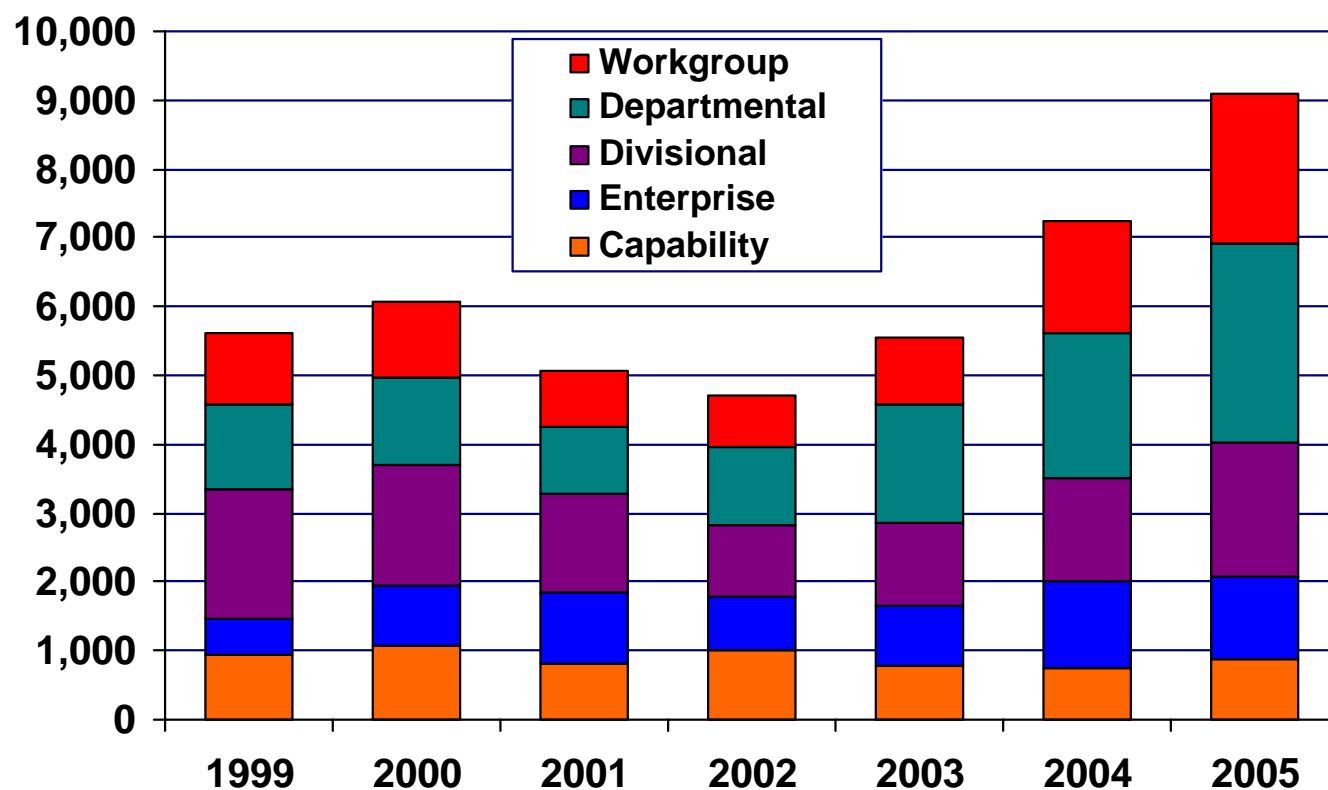
The New Realities

- **Major market growth over the last 4 years**
 - 94% growth since 2002
 - 23% growth in 2005 – Now **\$9.1 billion a year**
- **Clusters have been a disruptive force**
 - 1/3 of the market in 2004
 - Now close to 1/2 of the market
 - Caused a growth revolution, not a decline
- **Capability market transition continues**
 - Down -13% since 2002
- **Strong growth at the lower end of the market**
 - Workgroup up 200%, Departmental up 155%, Divisional up 84% since 2002
- **Bio-Sciences & government markets are growth areas**

Slide courtesy of Earl Joseph Jr. , IDC



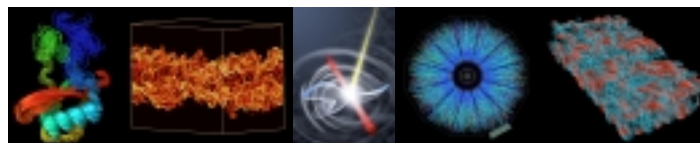
HPC Revenue by Competitive Segment (\$K)

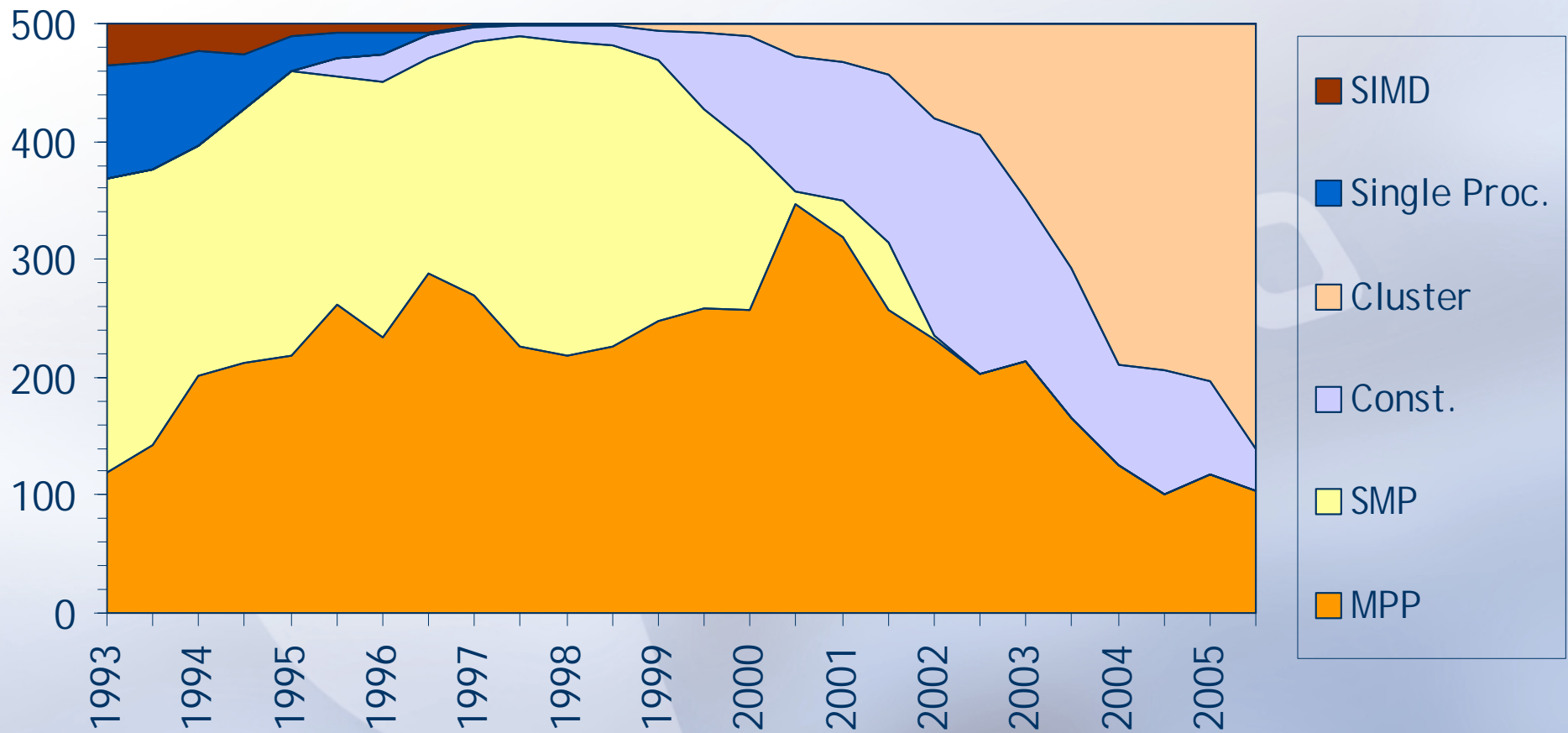


Revenue growth in 2005 favored the low-end:

- Capability declined -2.0%
- Enterprise declined -3%
- Divisional grew 30%
- Departmental grew 36%
- Workgroup grew 33%

Slide courtesy of Earl Joseph Jr. , IDC

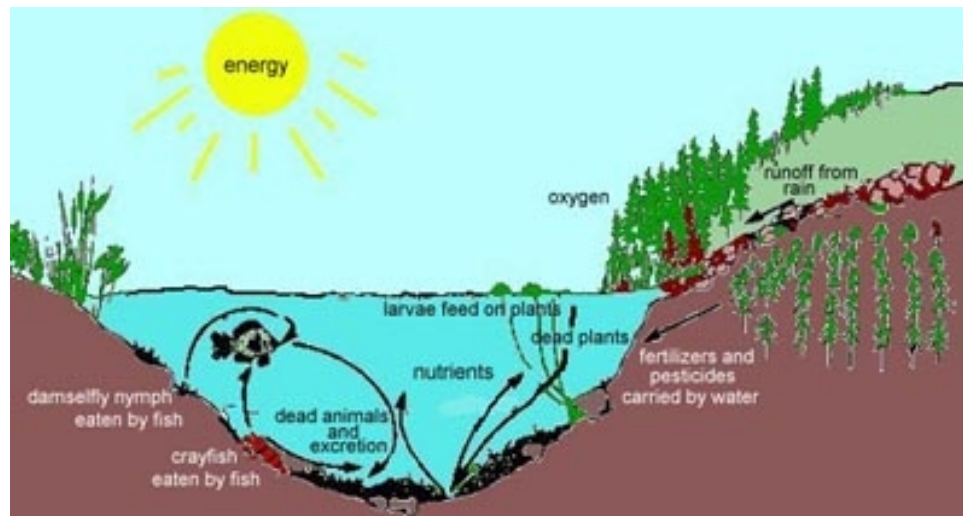




Ecosystem for HPC

From the NRC Report on “The Future of Supercomputing”:

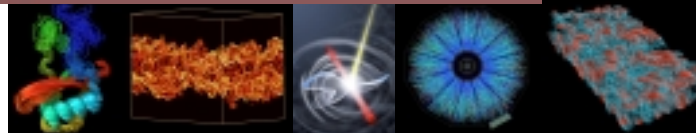
- Platforms, software, institutions, applications, and people who solve supercomputing applications can be thought of collectively as an ecosystem
- Research investment in HPC should be informed by the ecosystem point of view - progress must come on a broad front of interrelated technologies, rather than in the form of individual breakthroughs.



Pond ecosystem image from

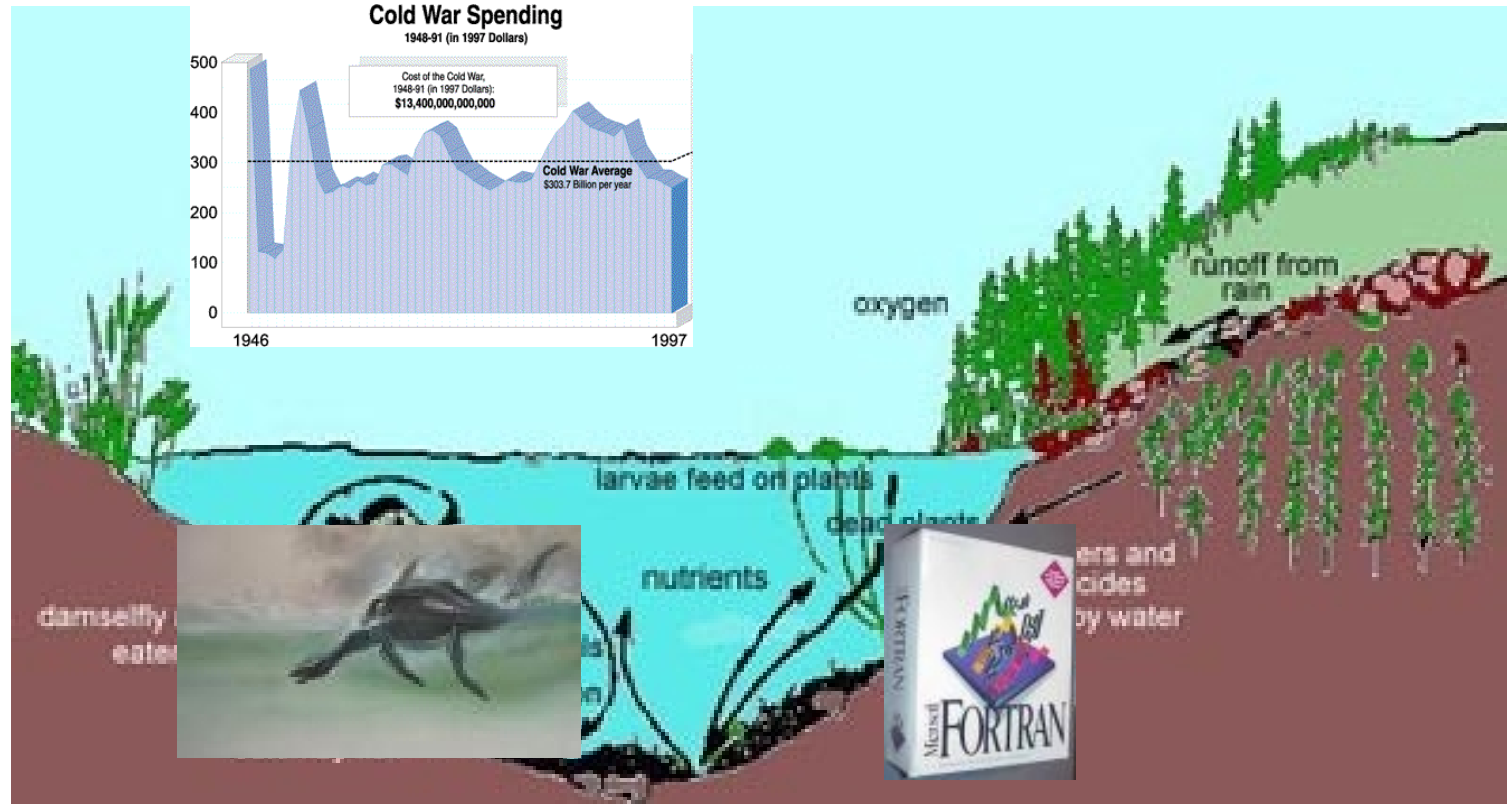
<http://www.tpwd.state.tx.us/expltx/e>

[ft/txwild/pond.htm](http://www.tpwd.state.tx.us/expltx/e/ft/txwild/pond.htm)



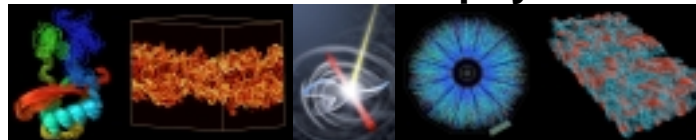
Supercomputing Ecosystem (1988)

Cold War and Big Oil spending in the 1980s



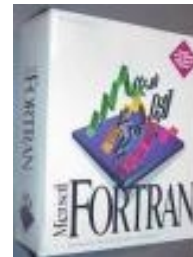
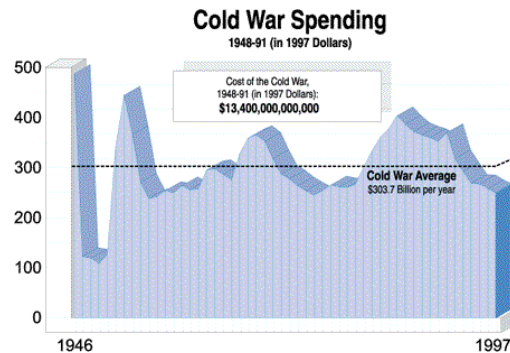
Powerful Vector Supercomputers

20 years of Fortran applications base in
physics codes and third party apps

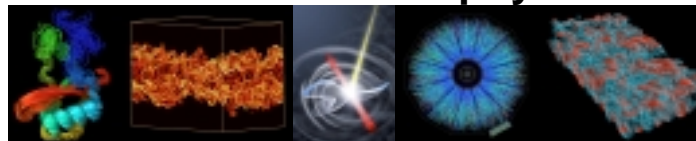


Supercomputing Ecosystem (until about 1988)

Cold War and Big Oil spending in the 1980s



Powerful Vector Supercomputers

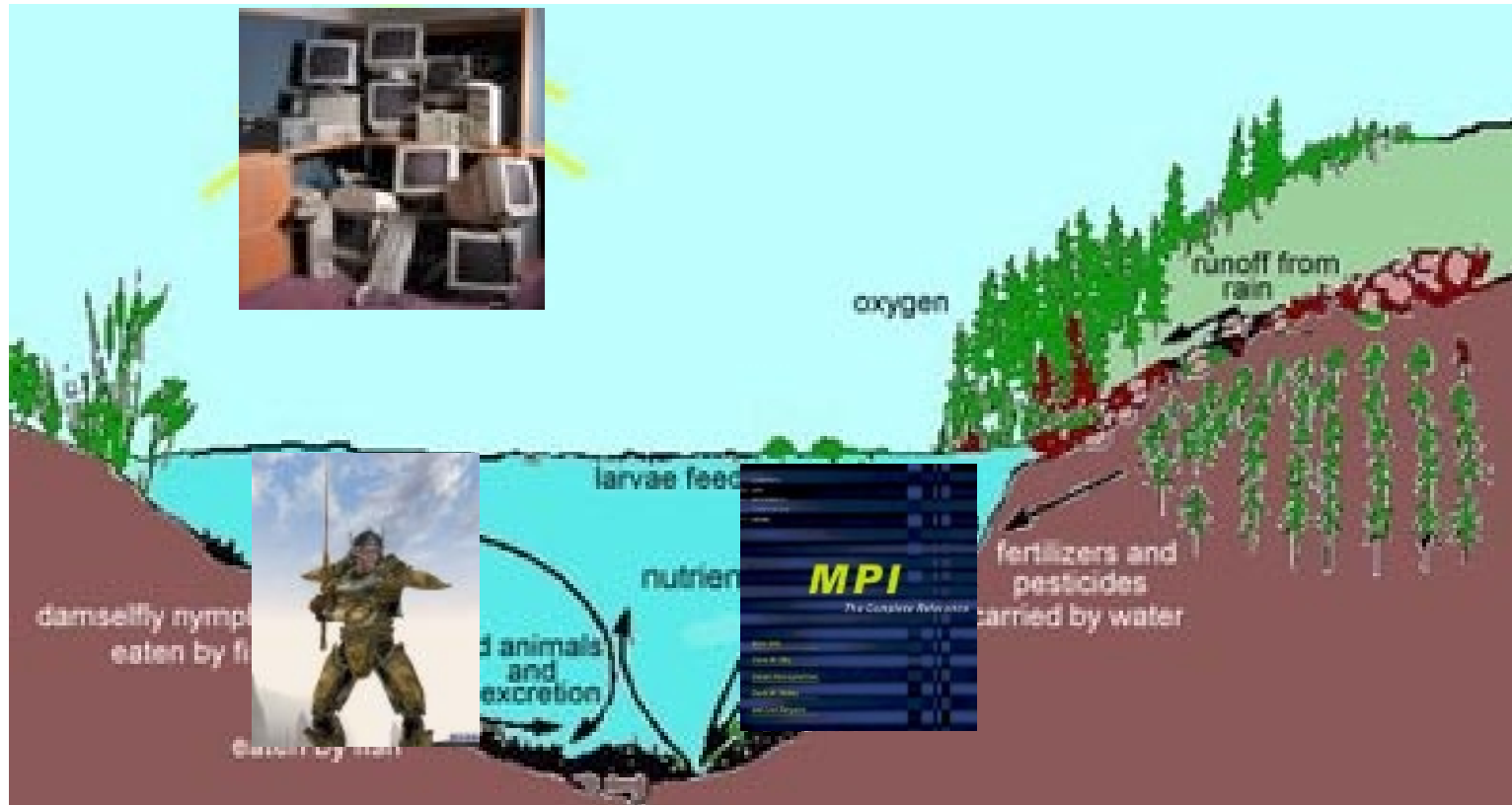


20 years of Fortran applications base in physics codes and third party apps



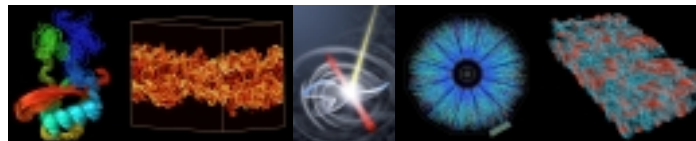
Supercomputing Ecosystem (2006)

Commercial Off The Shelf technology (COTS)



“Clusters”

12 years of legacy MPI applications base



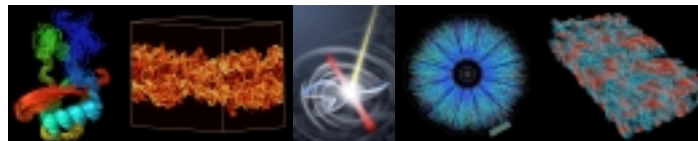
Supercomputing Ecosystem (2006)

Commercial Off The Shelf technology (COTS)



“Clusters”

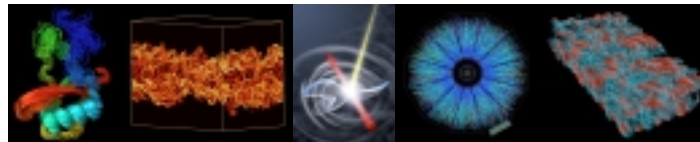
12 years of legacy MPI applications base



Observations on the 2006 Ecosystem

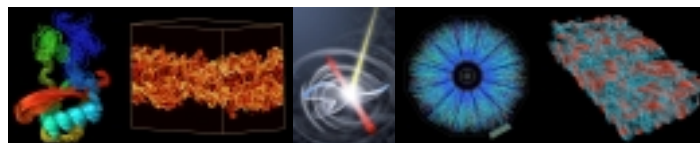
- **It works very well**
 - As long as you are content with <1000 processors
- **It is rapidly expanding**
 - IDC data on technical computing market growth
- **It is very stable and thus stifles innovation**
 - attempts of re-introducing “old species” not successful (X1)
 - attempts of introducing new species failed (mutation of Blue Gene/Cyclops 1999 to BG/L 2005)

The economic and applications success of clusters threatens innovations at the high end



Overview

- History and Future of Petaflops Computing
- HPC in 2006: “It was the best of times, it was the worst of times ...”
- “A Petaflops before its Time”
- The power problem
- The scaling problem
- What’s next?



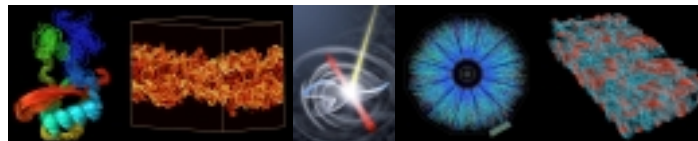
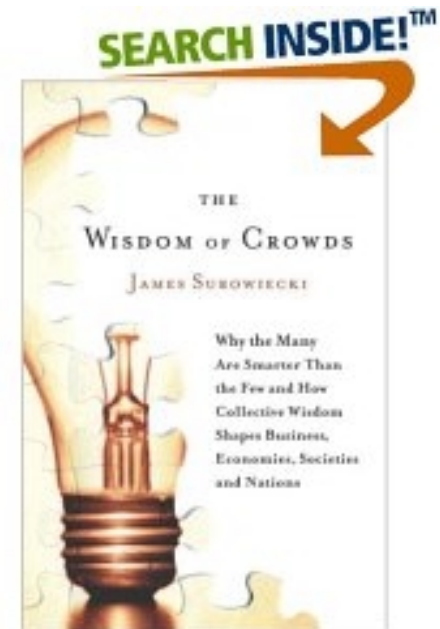
The Dagstuhl Experiment

**Based on the “Wisdom of Crowds”
by James Surowiecki.**

**In February 2006 I asked about 40
experts in HPC five simple questions
about Petascale computing.**

**All were participants of the (by invitation
only) Schloss Dagstuhl workshop on
Petascale Algorithms and Application**

**I collected the answers and report the
average as the collective opinion of
this workshop.**

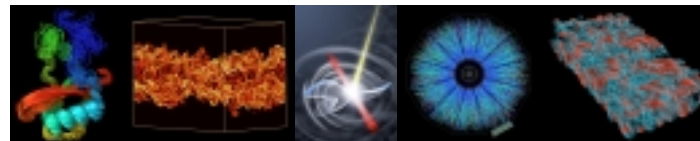


Q1: When will the first Petaflops machine be installed?

More precisely ...

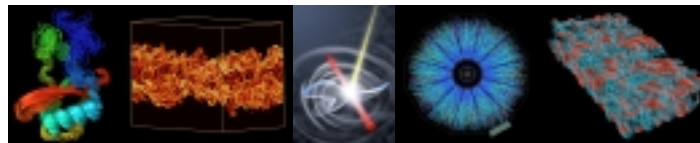
Question 1:

What date will the first machine appear on the TOP500 list that has a LINPACK RMAX > 1 Pflops?



Hints for Question 1

- The TOP500 list appears twice a year in June and in November
- First systems
 - BG/L >100 Tflops in 6/2005
 - Earth Simulator > 10 Tflops in 6/2002
 - ASCI Red > 1 Tflops in 6/1997
 - Numerical Windtunnel > 100 Gflops in 11/1993

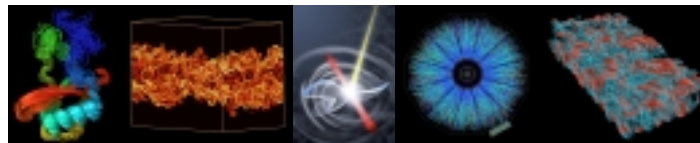


Q2: When will the first Petaflops application performance be obtained?

More precisely ...

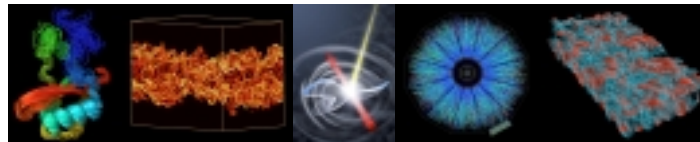
Question 2:

What date will the first Gordon Bell Prize be awarded to an application that performs in excess of 1 Pflops?



Hints for Question 2

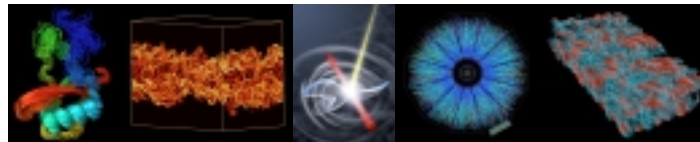
- **The Gordon Bell Prize is awarded annually at the SCxx conference**
- **Gordon Bell has made an endowment to ACM that will guarantee funding for the prize for at least another 25 years.**
- **First applications to reach a level**
 - Molecular Dynamics > 100 Tflops in 2005 (BG/L)
 - Climate Simulation > 10 Tflops in 2002 (ES)
 - Material Science > 1 Tflops in 1998 (T3E)
 - Structures > 100 Gflops in 1994 (Paragon)
- **Answer to Q2 should probably not be an earlier date than answer to Q1.**



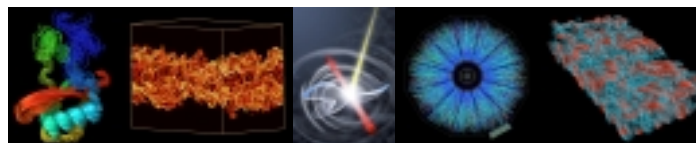
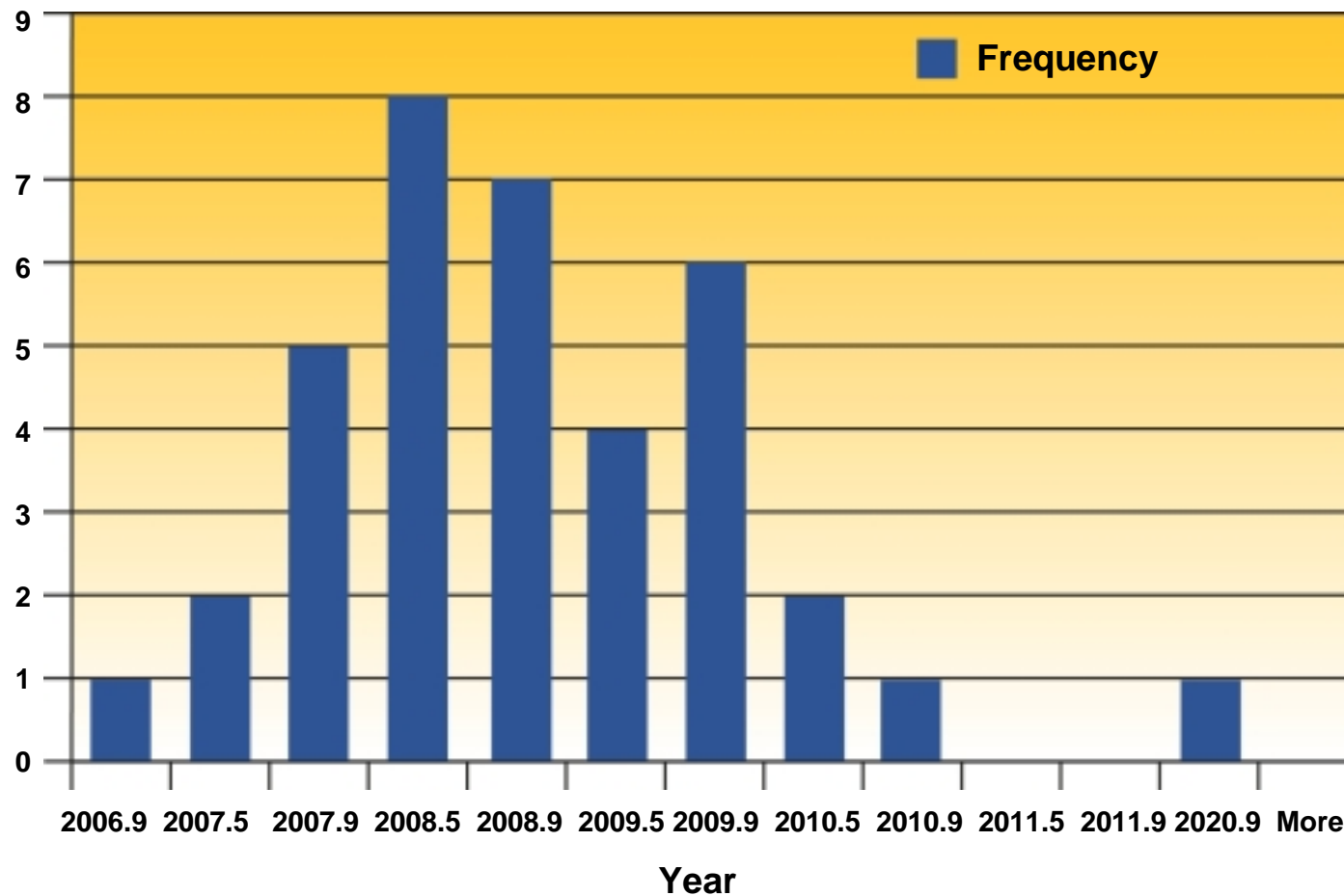
Bonus Question

Where will the first Petaflops system be installed?

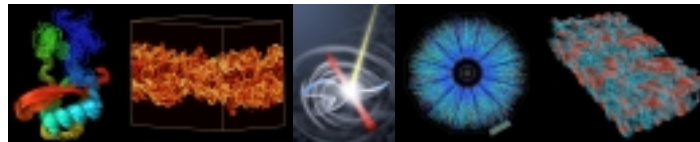
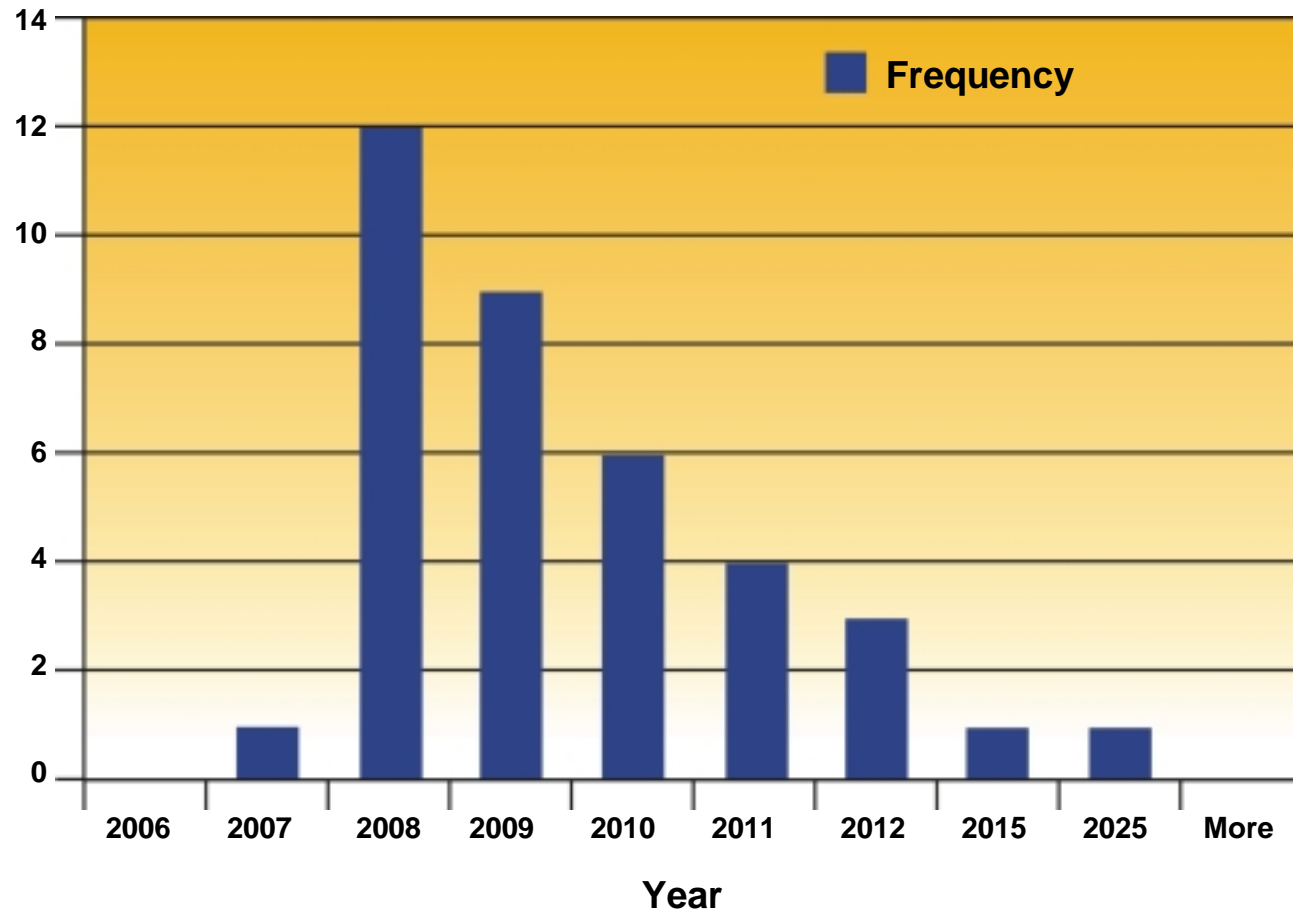
- A. China**
- B. EU (optional: which country?)**
- C. Japan**
- D. USA (optional: which state?)**
- E. other**



First Petaflops System on TOP500

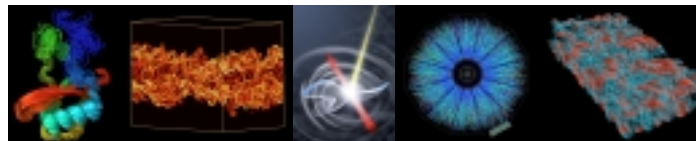
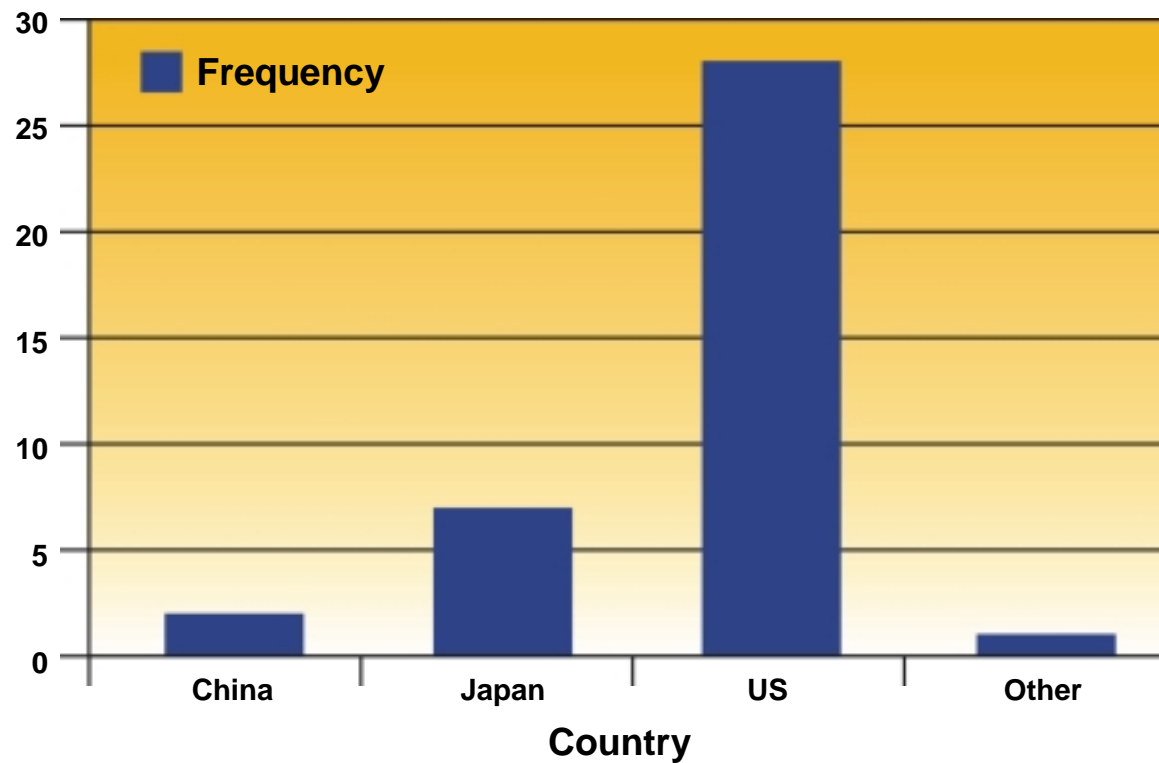


First Petascale Application



Country of Installation

Country of Installation of First Petascale System



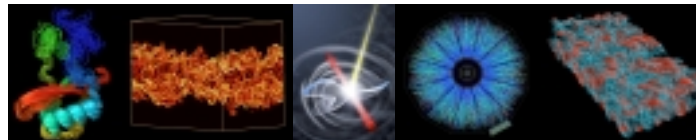
Increasing Blue Gene Impact

- **SC 2005 Gordon Bell Award, 101.7 TFs on real materials science simulation**
 - Recently exceeding 200 TFs sustained
- **Sweep of the all four HPC Challenge class 1 benchmarks**
 - G-HPL (259 Tflop/s), G-RandomAccess (35 GUPS) EP-STREAM (160 TB/s) and G-FFT (2.3 Tflop/s)
- **Over 80 large-scale applications ported and running on BG/L**



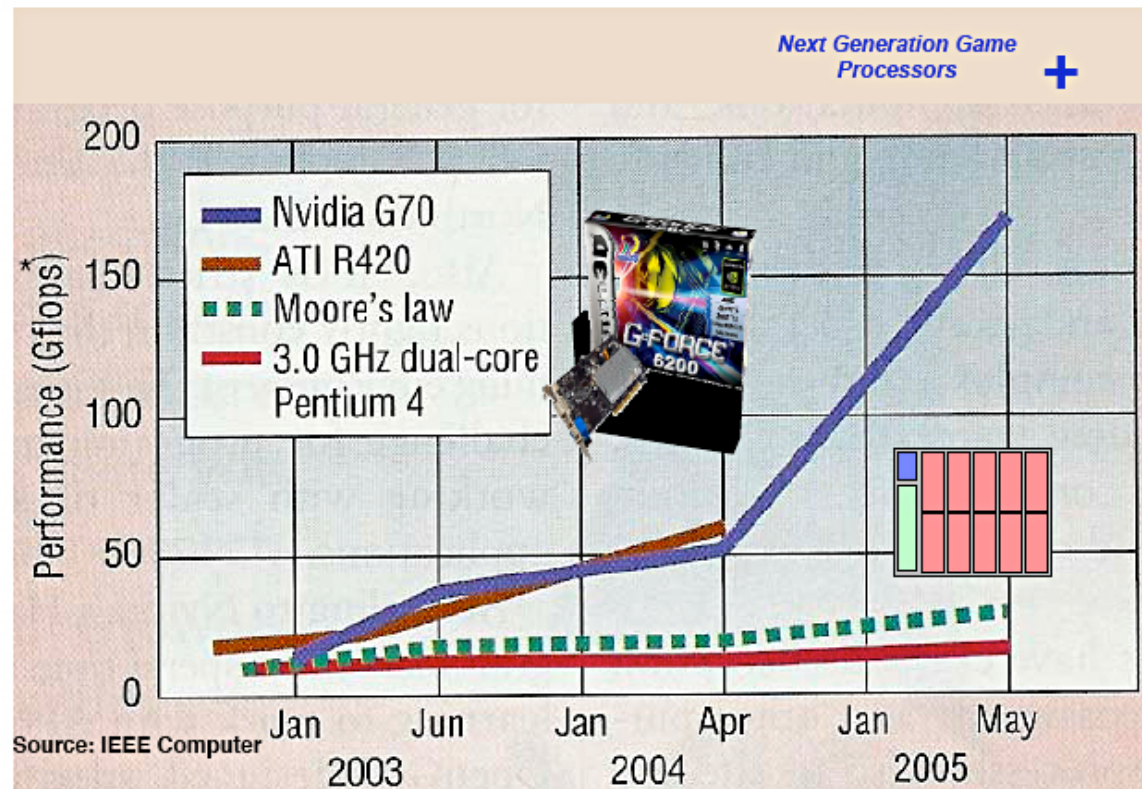
27.6 kW power
consumption per rack
(max)
7 kW power
consumption (idle)

Slide adapted from Rick Stevens, ANL



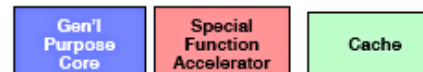
... and Increasing Game Processor Impact

- GPUs & Game Processor Architectures Are an Excellent Match for Game Applications
- Performance Has Been Growing Faster Than Moore's Law !

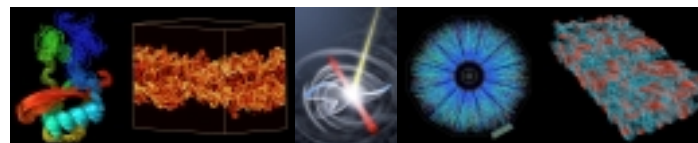


*Single Precision

• Typical Street Price for High End PC Graphics Card: **300 – 400 US\$**



Source: Randy Moulic, IBM

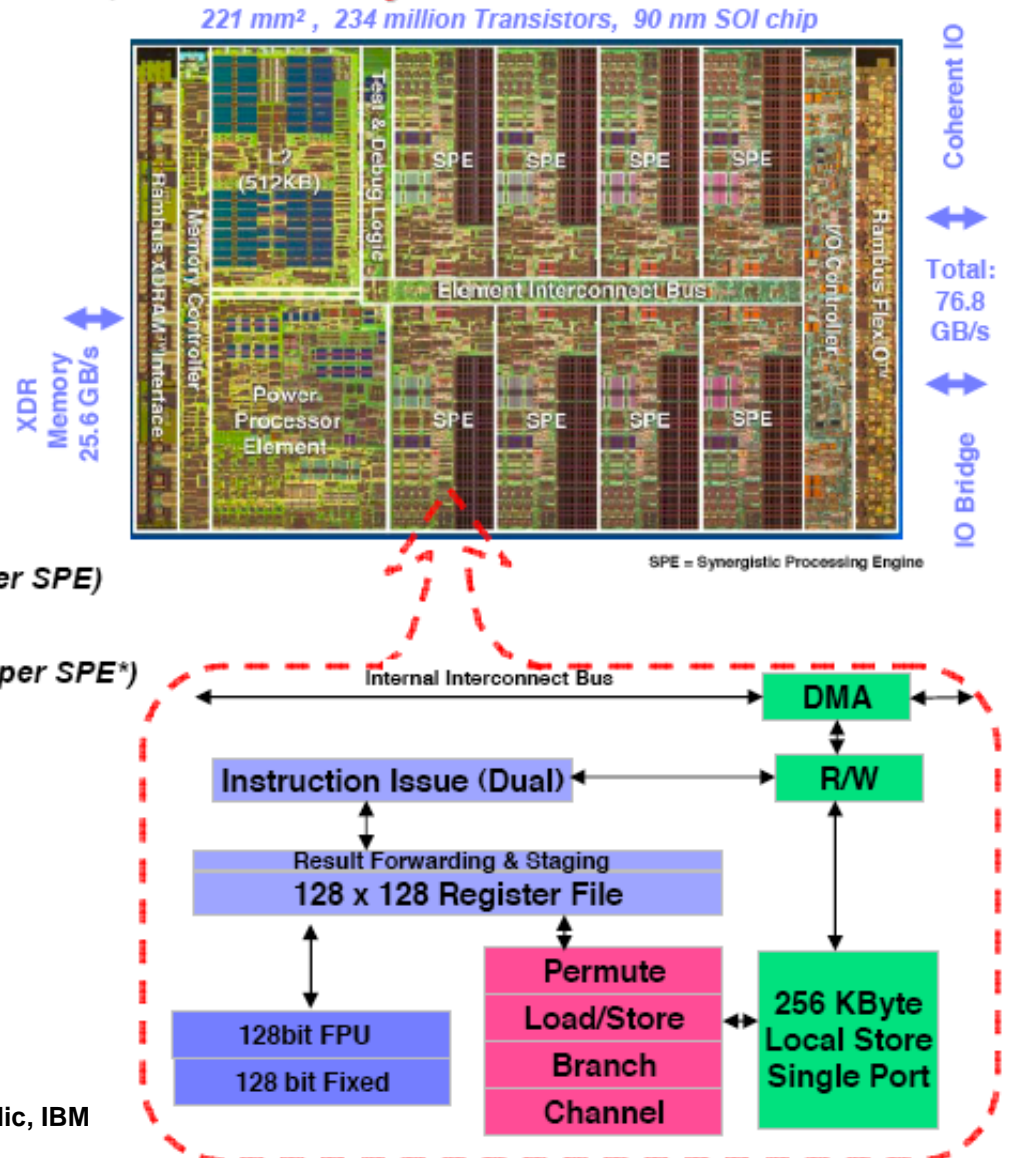


IBM Cell

Supercomputing Capability for High Volume, Consumer Systems

- Multi-core, multi-thread, "cluster-on-a-chip"
 - 64bit PowerPC Control Processor
 - + 8 Tightly integrated accelerators (SPE)
 - 128 bit SIMD/Vector, MAC
 - 256KB Embedded Memory
 - Integrated I/O and memory interfaces
- High Performance
 - 3.2 GHz clock frequency
 - 205 GFLOPs/s peak, single precision
(dual issue, in-order execution, 25.6 GFLOPS per SPE)
 - ~20 GFLOPs/S peak, double precision
(2 DP instructions every 7 cycles, 1.83 GLOPS per SPE*)
 - 205 GB/s internal interconnect bandwidth
 - ~100 GB/s BW for memory, external IO
- Linux OS
 - Simultaneous multiple OS support
 - + Real-time support

Source: Randy Moulic, IBM



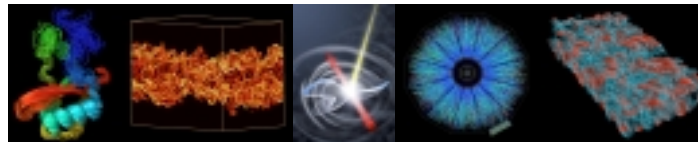
The Allure of Game Processors

The Potential of the Cell Processor for Scientific Computing

Samuel Williams, John Shalf, Leonid Oliker
Shoaib Kamil, Parry Husbands, Katherine Yelick
Computational Research Division
Lawrence Berkeley National Laboratory
Berkeley, CA 94720

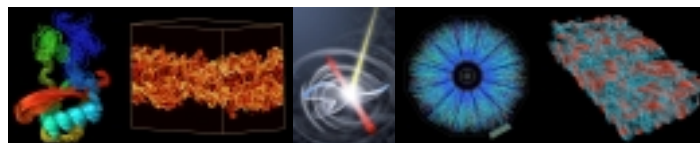
{swwilliams,jshalf,loliker,sakamil,prjhusbands,kayelick}@lbl.gov

- 30,000 downloads in one week
- Most accessed article on HPCWire **EVER!**



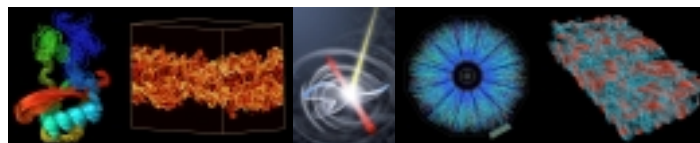
A Petaflops before its Time

- Even among experts there is an undue optimism about how close we are to “Petascale” computing
- In 11/2008 there will be a (Linpack Rmax) Petaflops computer on the TOP500 list
- Most likely it will be a BG/P, possibly a “souped up” Cell
- It will create an unwarranted sense of accomplishment
- It will distract from the development of real production Petaflops systems (e.g. XT, Power)



Overview

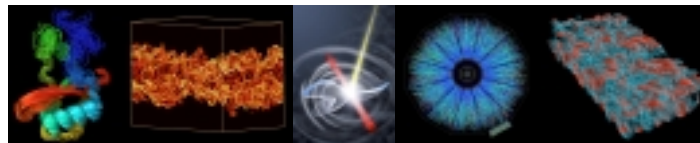
- History and Future of Petaflops Computing
- HPC in 2006: “It was the best of times, it was the worst of times ...”
- “A Petaflops before its Time”
- The power problem
- The scaling problem
- What’s next?



NERSC Estimate ...

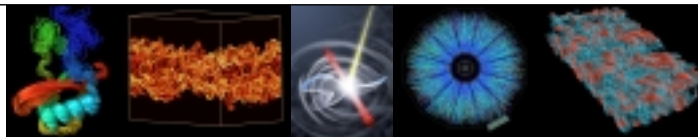
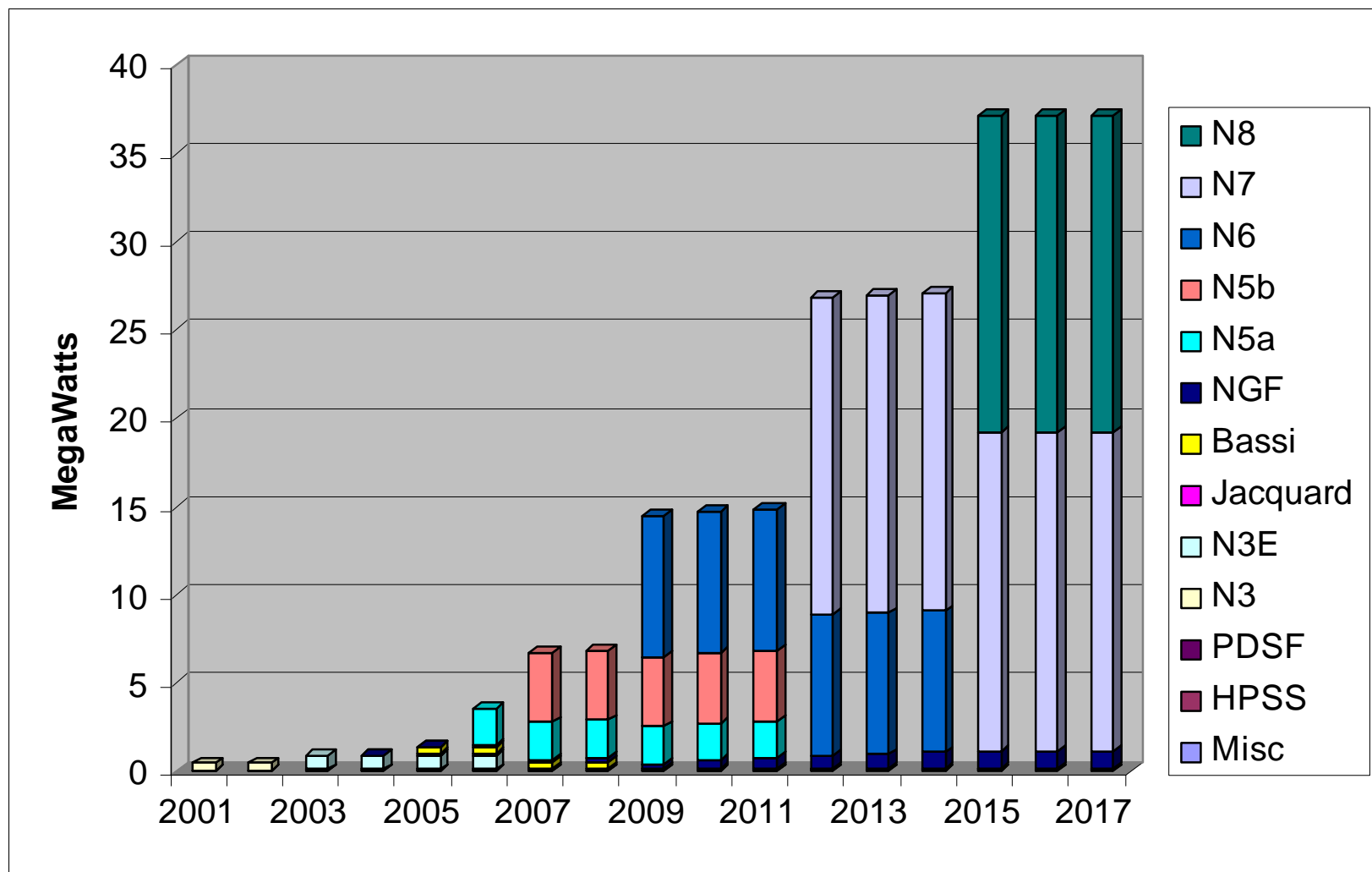
... for a sustained Petaflops system (on multiple applications) in 2010

- **20 MW**
- **16,000 square feet**
- **\$12M/year electricity cost**

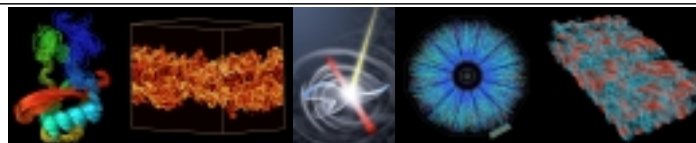
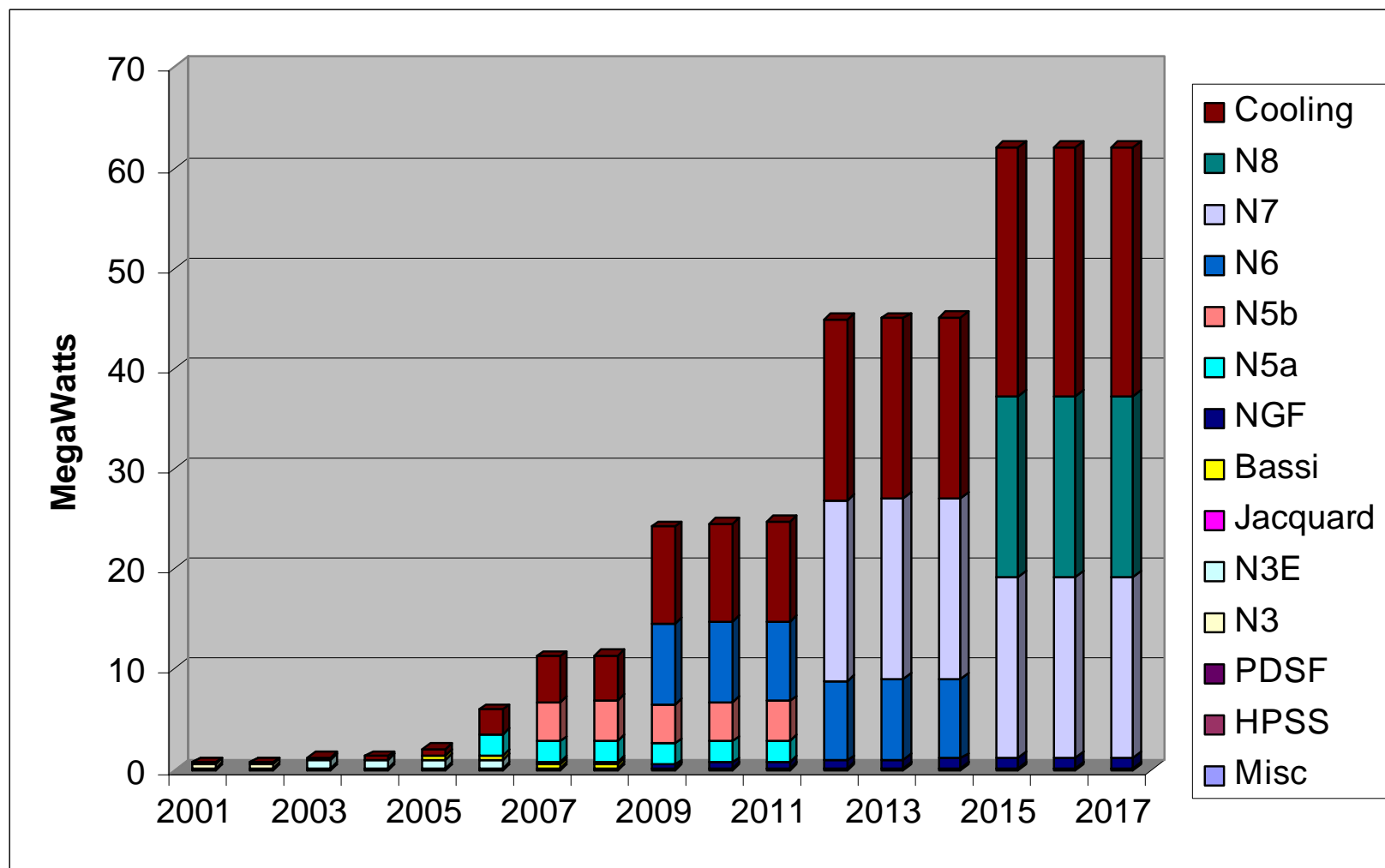


NERSC Projections for Computer Systems Power

(Does not include cooling)

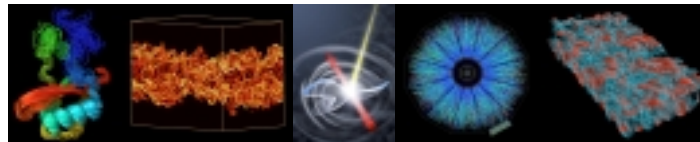


NERSC Projections for Computer Room Power System + Cooling

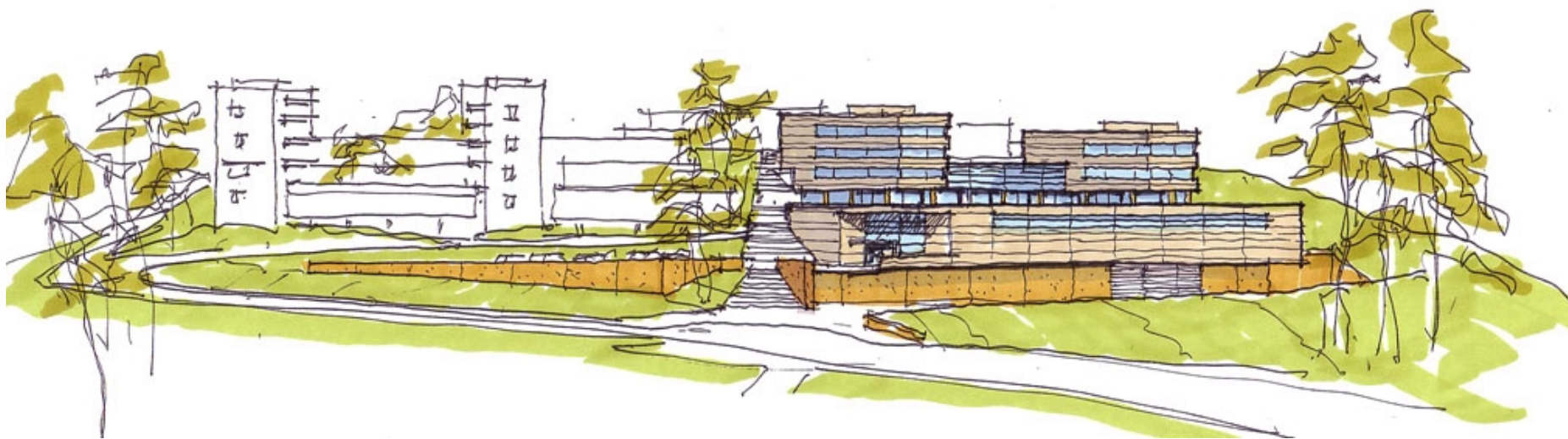


How Electrical and Space Costs Will Evolve

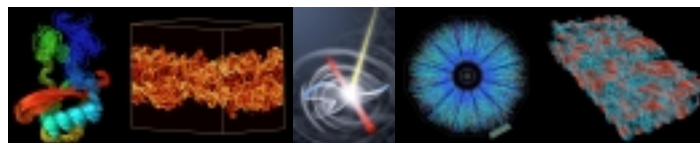
- **Price performance (Flop/s per \$) increasing faster than the facility needs (Flop/s per sf and Flop/s per watt)**
- **Liquid Cooling Tradeoffs:**
 - allows higher density, and some cooling efficiency
 - Less space but more electricity
- **Configurations impose constraints on layout**
 - More space use due to cable limitations and topology



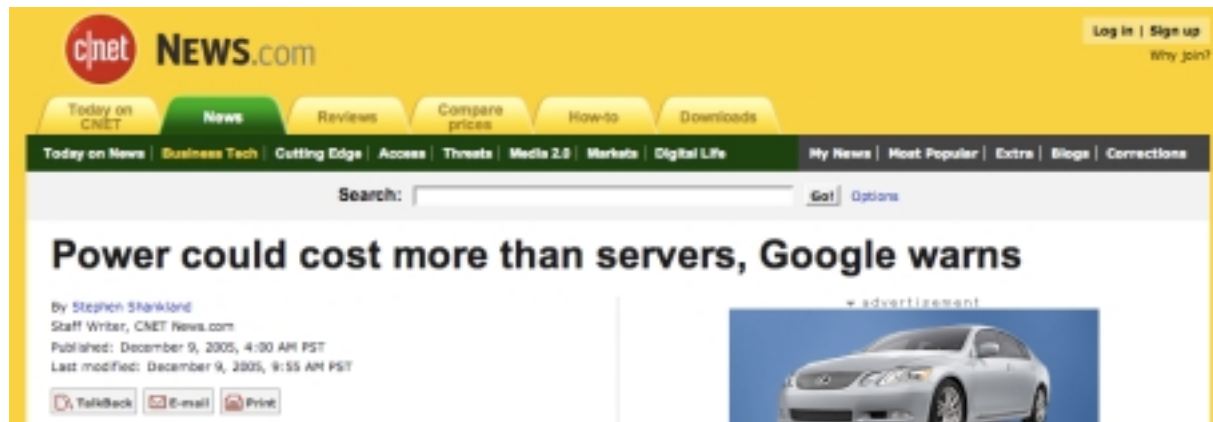
New Building in Berkeley



... but new \$90M buildings cannot be the industry solution



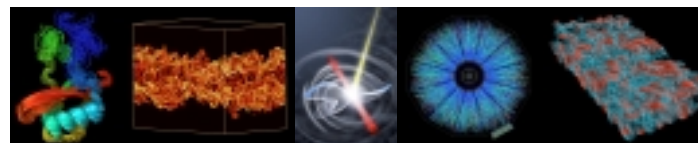
Power will be an Industry Wide Problem



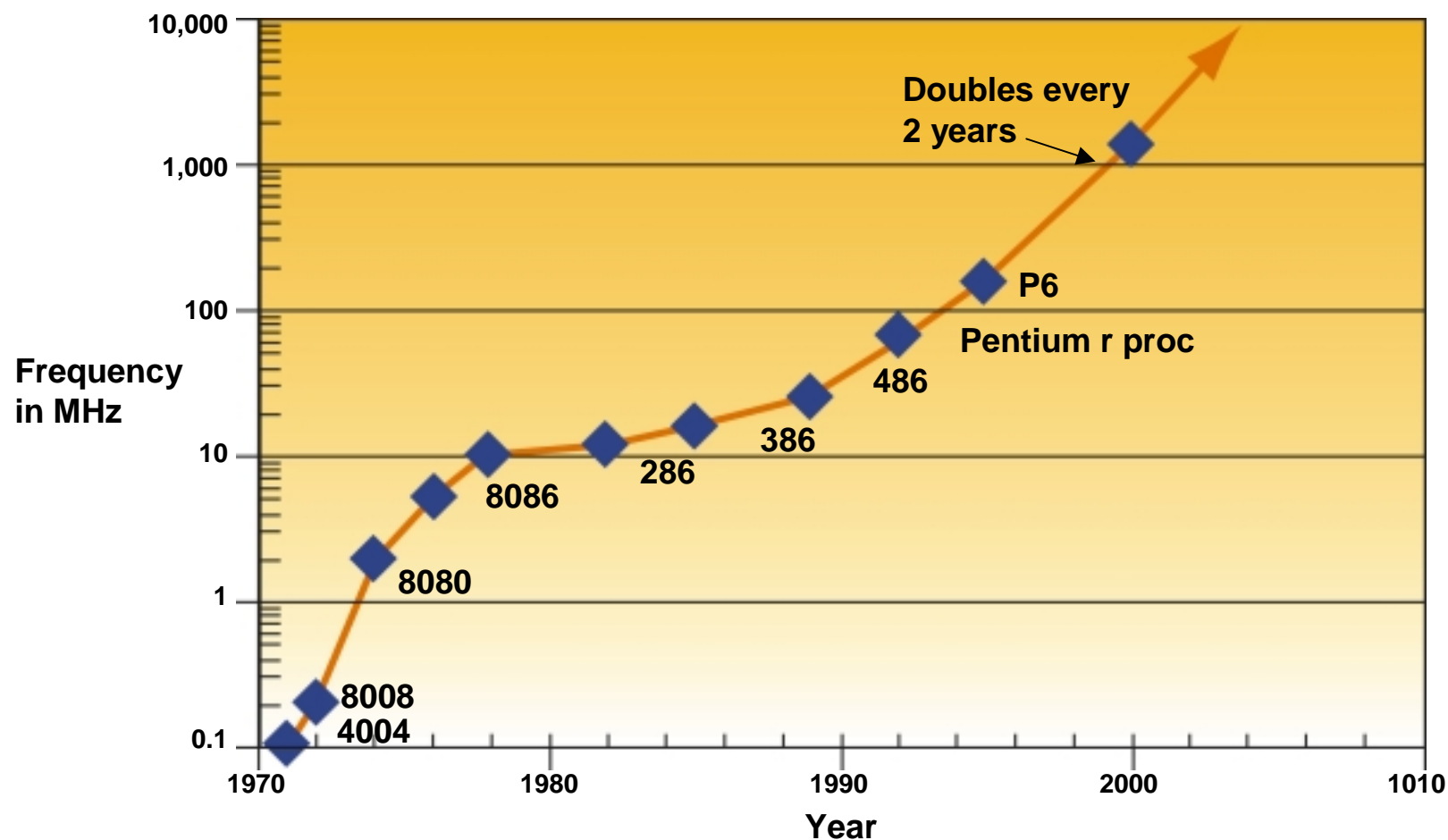
The New York Times **“Hiding in Plain Sight, Google Seeks More Power”,**
by John Markoff, June 14, 2006



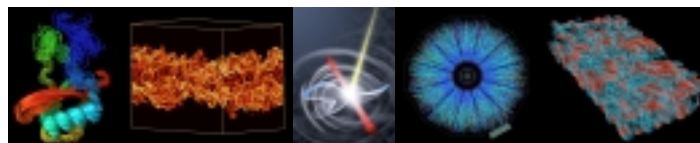
New Google Plant in The Dalles, Oregon,
from NYT, June 14, 2006



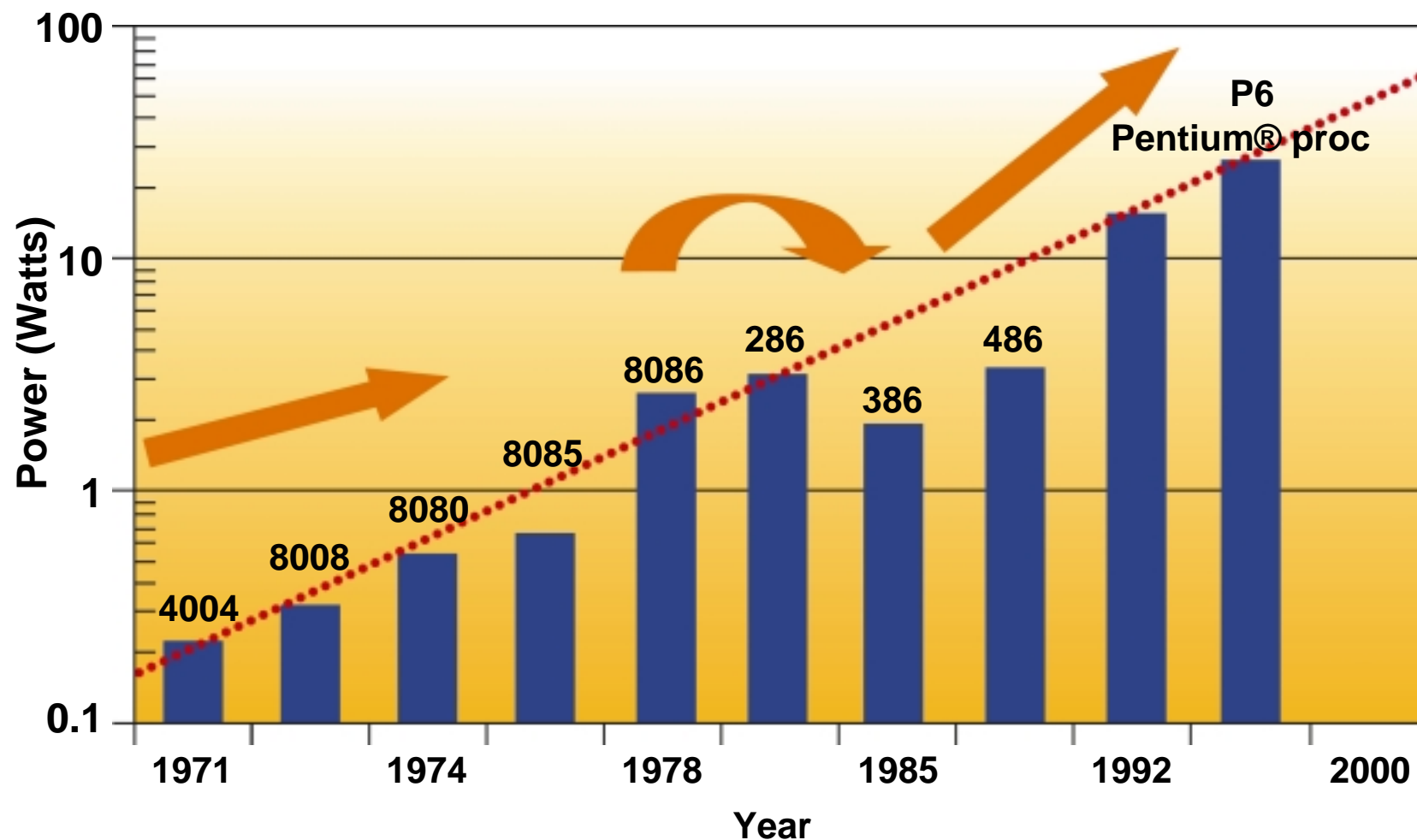
Intel Prediction of Microprocessor Frequency (ca. 2001)



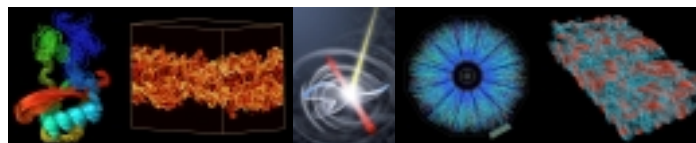
Adopted from a presentation by S. Borkar, Intel



Intel Prediction of Microprocessor Power Consumption (ca. 2001)



Adopted from a presentation by S. Borkar, Intel



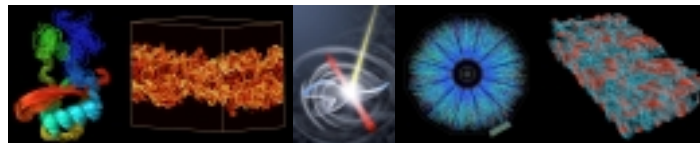
Learning from the History of Microprocessors

1999: nobody paid any attention to the issue of power

2001: first warnings about future power problems - Borkar speaks about heat dissipation equivalent to a rocket nozzle

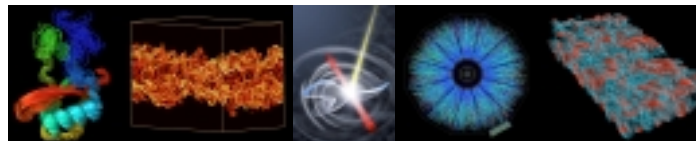
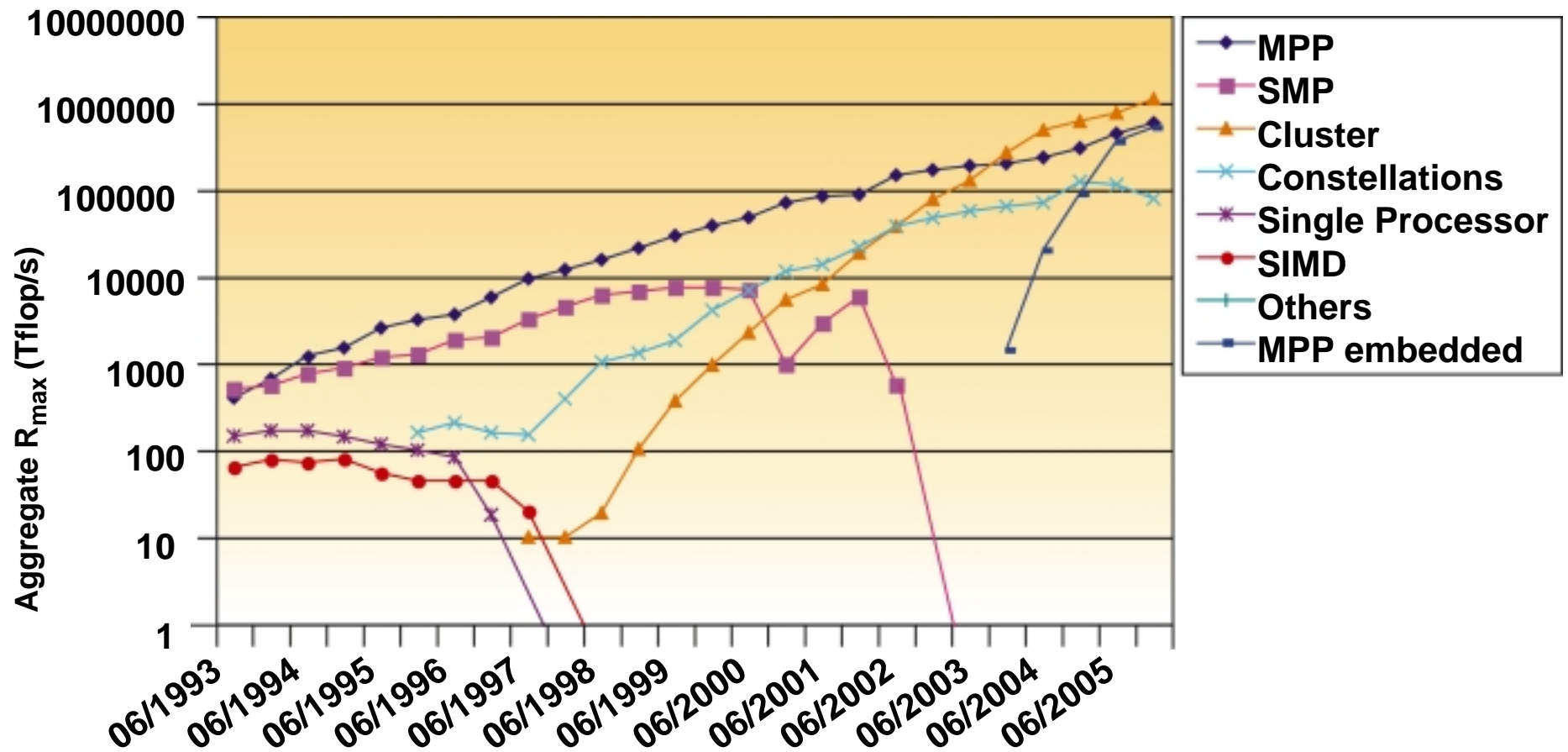
2004: Intel announces the end of the drive for more performance through increased clock rates

2006: Do we see any solutions for the systems power problem?



BG/L—the Rise of the Embedded Processor?

TOP 500 Performance by Architecture

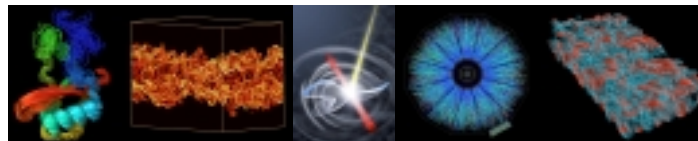


Future Scaling without Innovation

If we scale current peak performance numbers for various architectures and allowing system peak doubling every 18 months. **Trouble ahead**

	Projected Year	BlueGene/L	Earth Simulator	MareNostrum
250 TF	2005	1.0 MWatt	100 MWatt	5 MWatt
1 PF	2008	2.5 MWatt	200 MWatt	15 MWatt
10 PF	2013	25 MWatt	2000 MWatt	150 MWatt
100 PF	2020	250 MWatt	20,000 MWatt	1500 MWatt

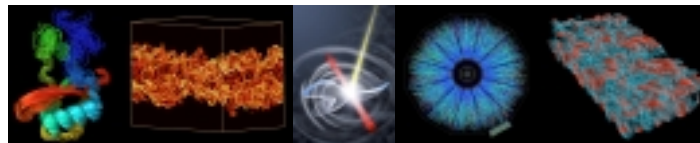
Slide adapted from Rick Stevens, ANL



Discussion About Power

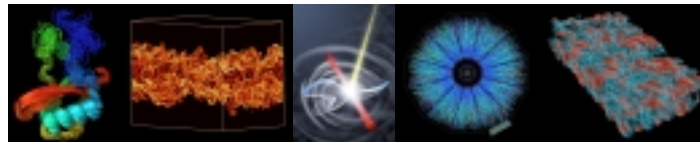
- In 2000 the “industry” expected a 10 GHz processor by about now - it did not happen
- In 2006 HPC experts are planning for 10-20 MW systems - will that really happen?
- Switching to embedded processor technology (a la BG) will possibly delay the problem by a few years
- Will we we run out of electricity before we run out of Moore’s Law?

Does it make sense to build systems that require the electric power equivalent of an aluminum smelter?



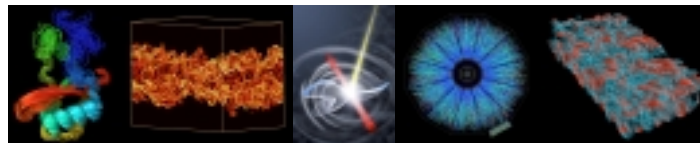
Overview

- History and Future of Petaflops Computing
- HPC in 2006: “It was the best of times, it was the worst of times ...”
- “A Petaflops before its Time”
- The power problem
- The scaling problem
- What’s next?

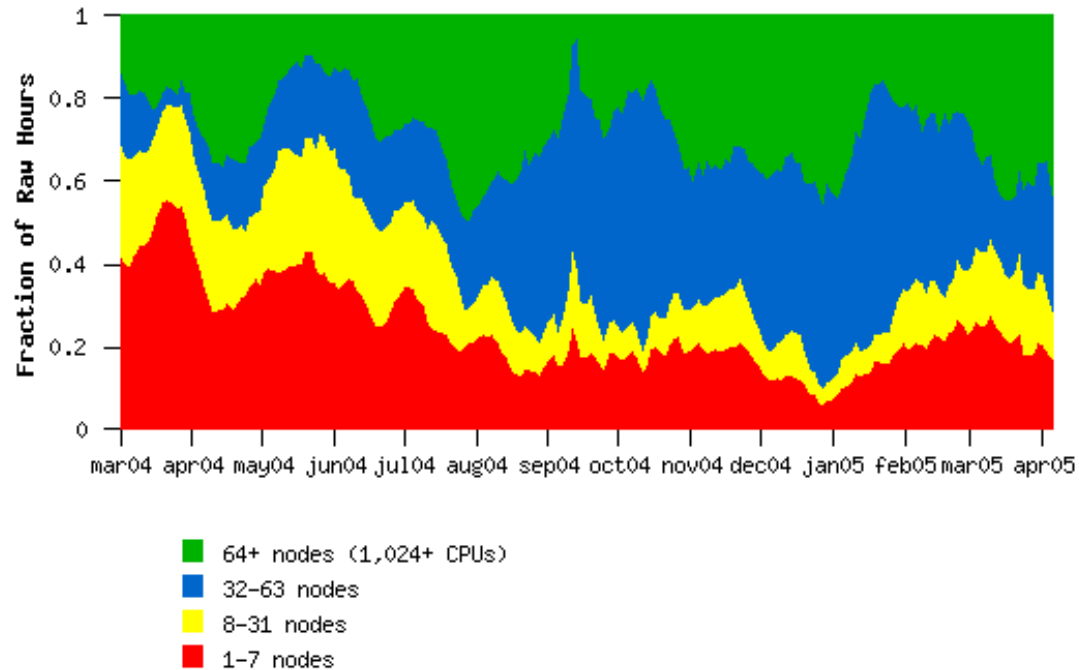


Scaling to Petaflops

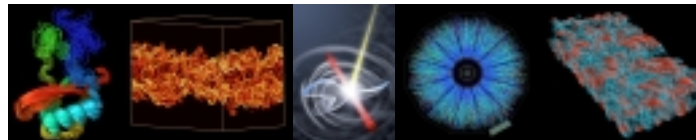
- **NERSC estimates that a sustained Petaflops system (on multiple applications) in 2010 will have 150,000 - 500,000 (multi-core) processors**
- **Today almost no applications and system software is ready to scale to that level**



Application Status in 2005

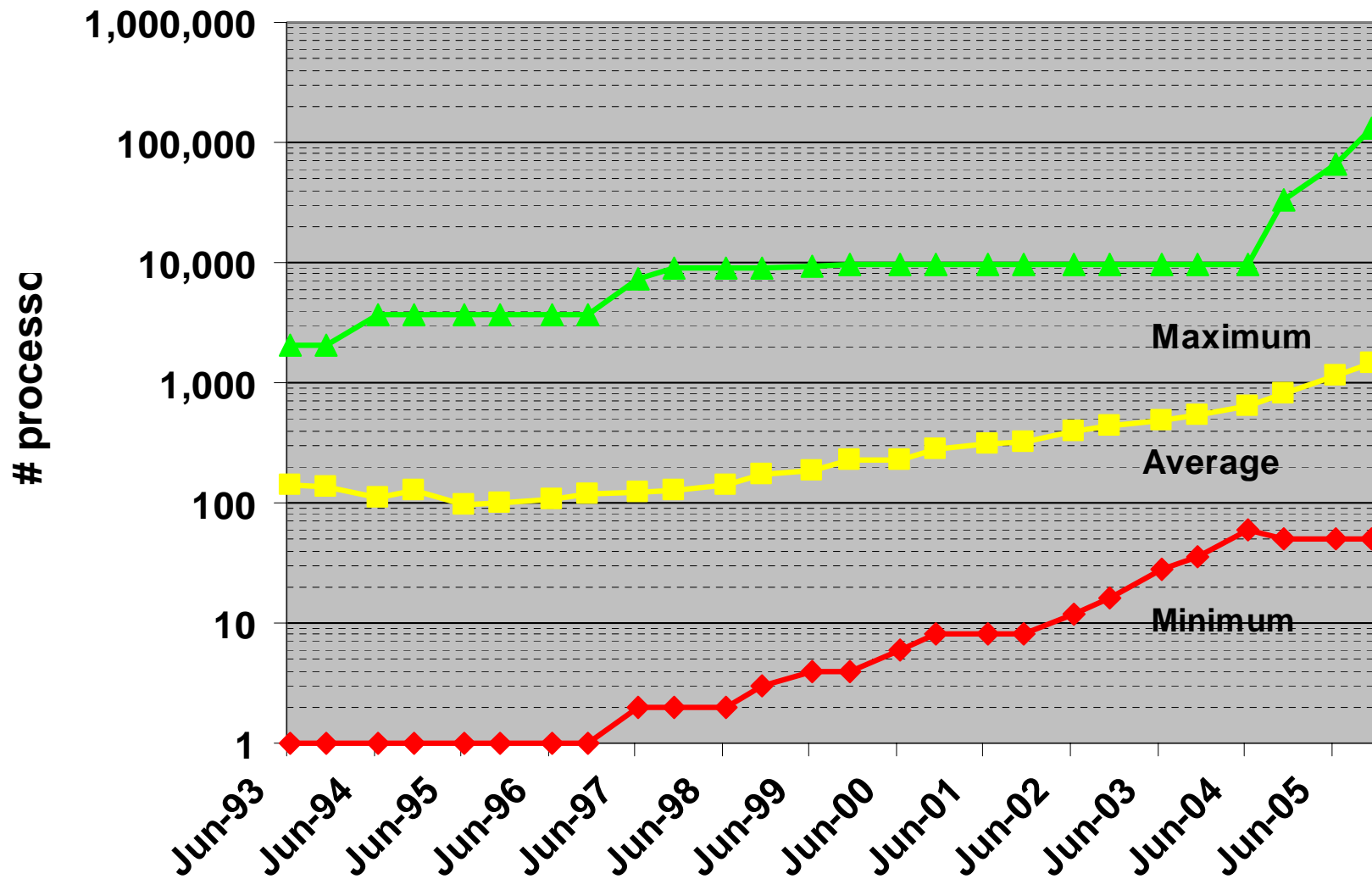


- A few Teraflop/s sustained performance
- Scaled to 512 - 1024 processors

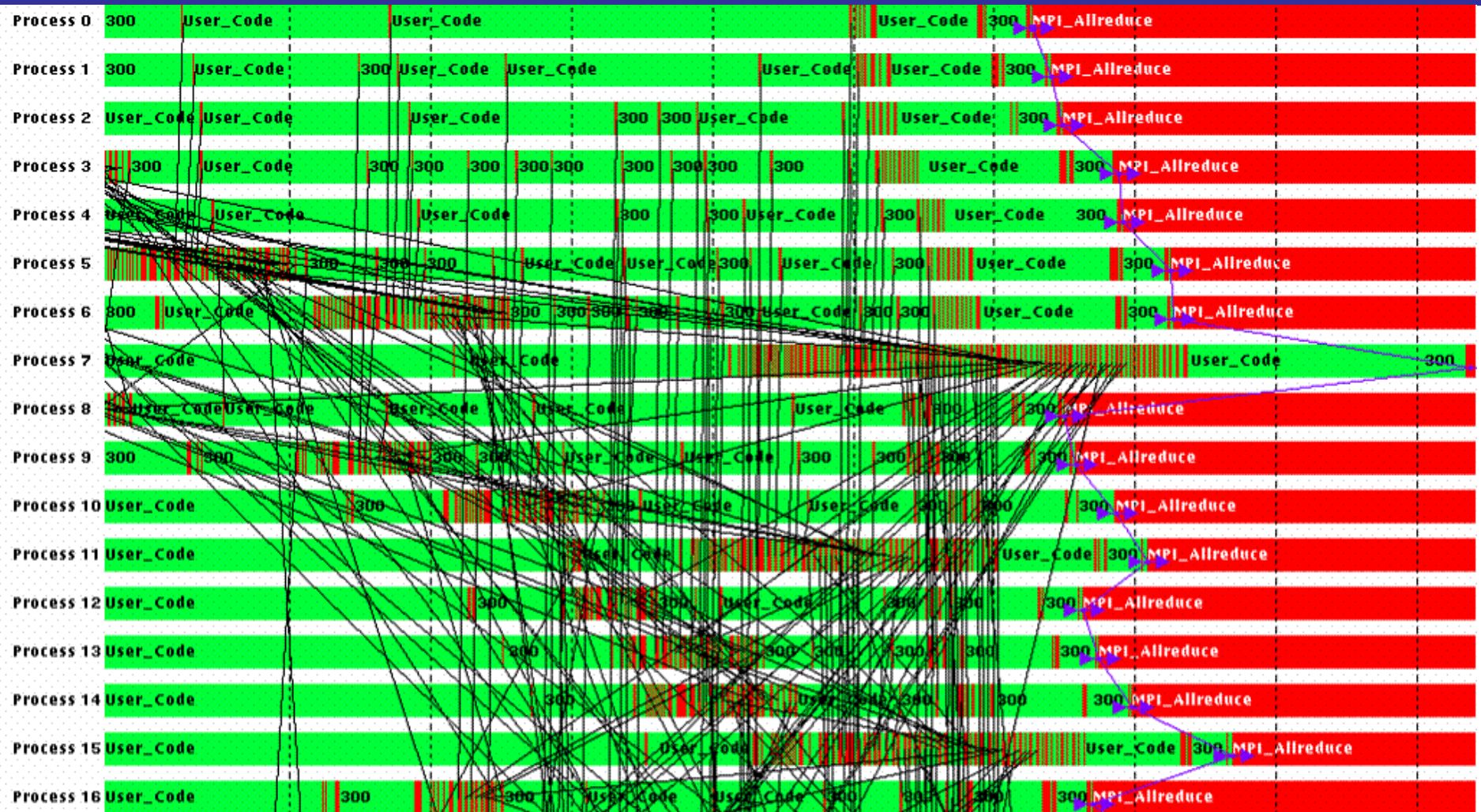


Parallelism has Stagnated for a Decade

Number of processors in the parallel system in the TOP500

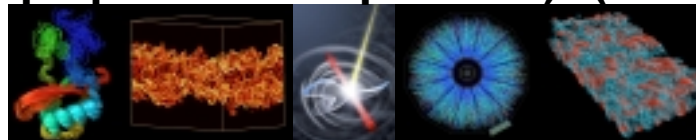


16 Way for 4 seconds

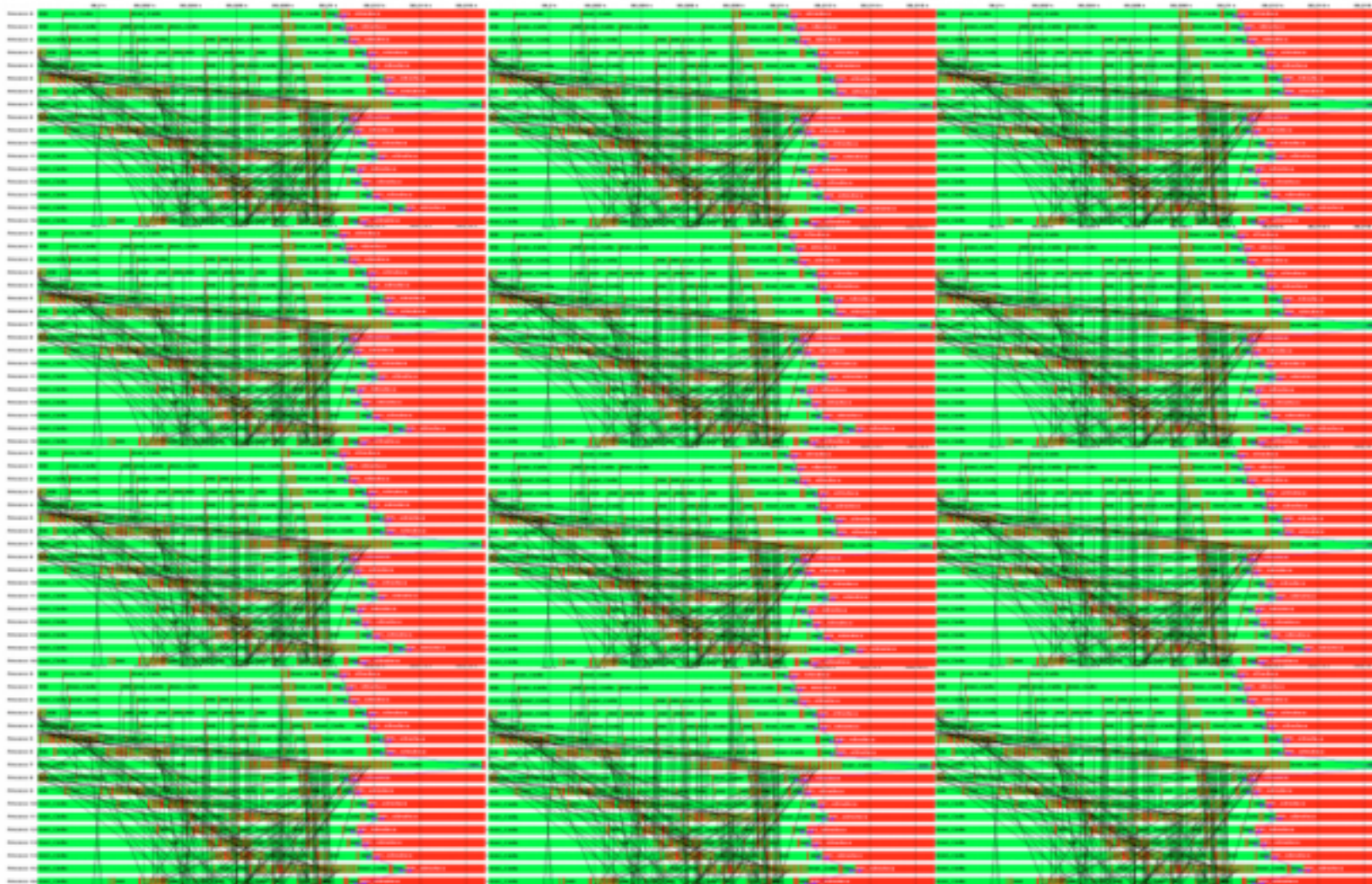


(About 20 timestamps per second per task) * (1...4 contextual variables)

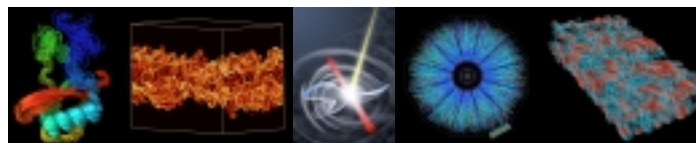
Slides by David
Skinner, NERSC



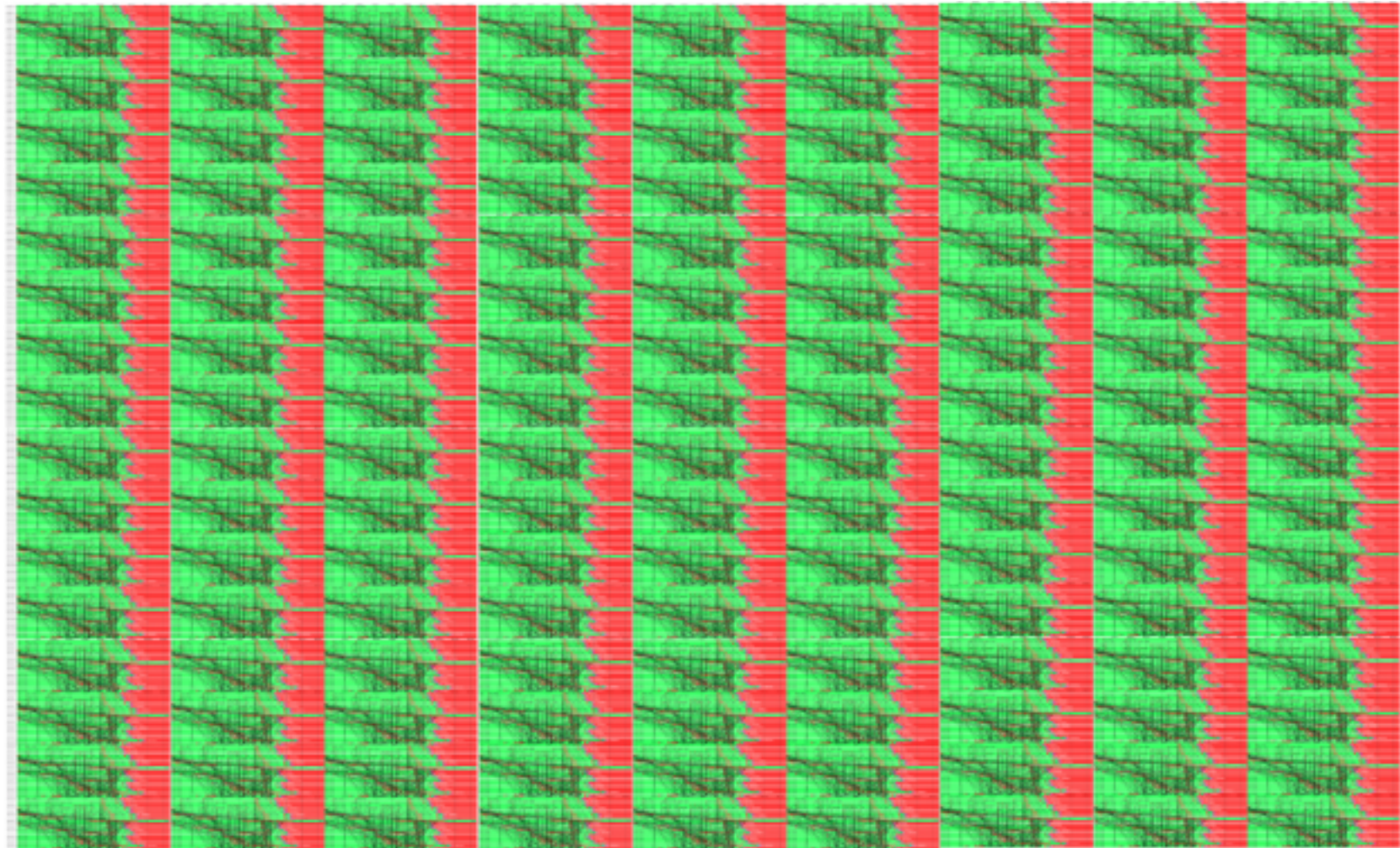
64 way for 12 seconds



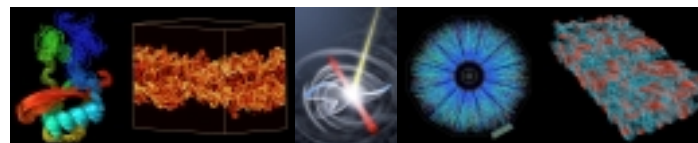
Slides by David
Skinner, NERSC



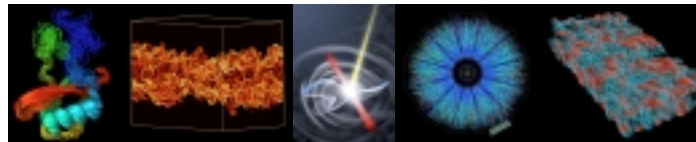
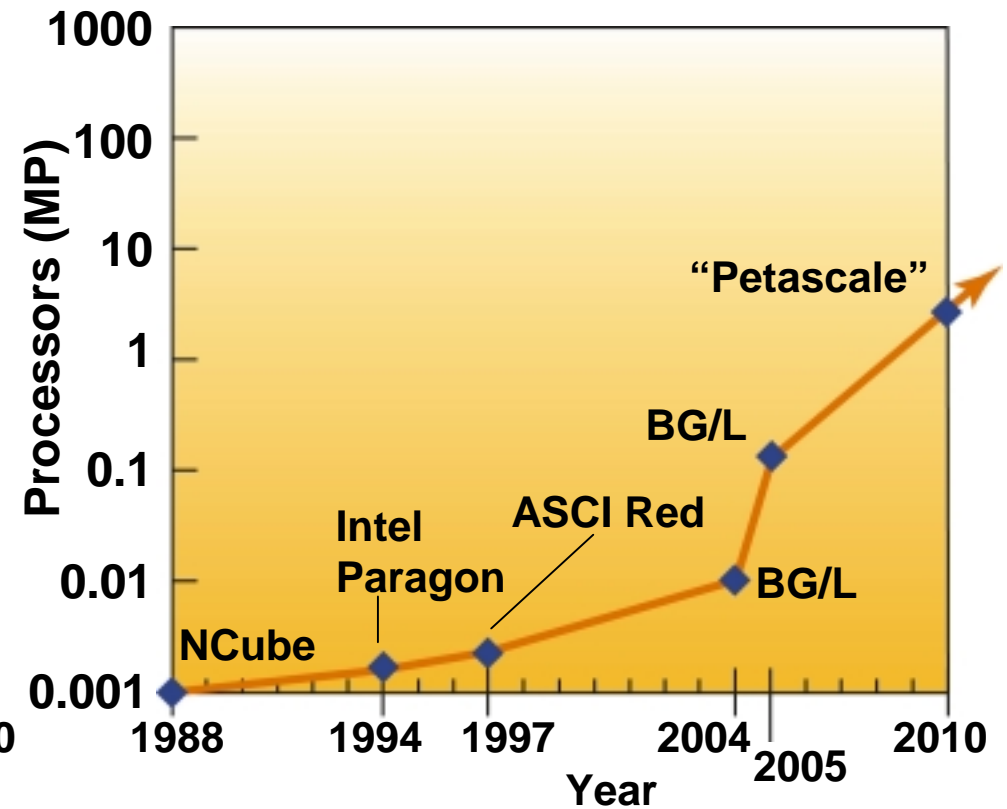
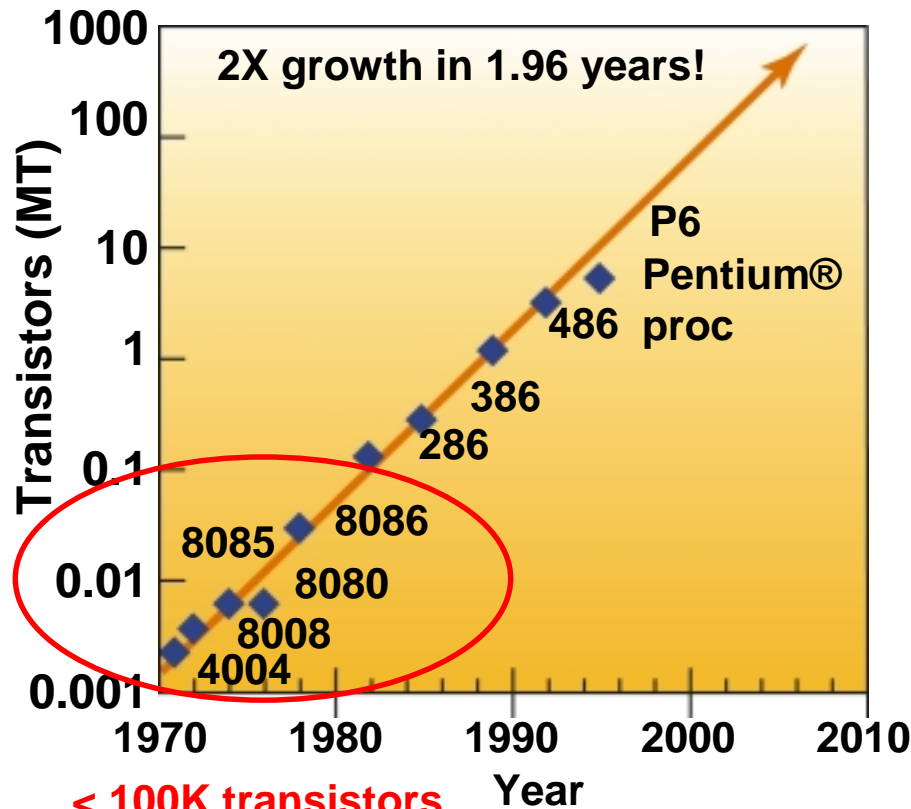
256 Way for 36 Seconds



Slides by David
Skinner, NERSC



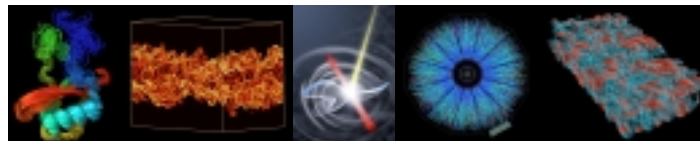
Historical Reference: Transistor Count



“The Processor is the new Transistor”

(David Patterson)

- **NERSC’s flagship computing system, Seaborg, contains as many processors as there are transistors in the original Intel 8080a implementation (6,000 transistors vs 6,000 processors)**
- **BG/L at LLNL contains as many processors as there are transistors in the MC68000 (manufactured in 1980, the MC68000L was a 32-bit processor and contained 68,000 transistors).**
- **With 1.5M processors, BG/Q likely to have more processors than there are logic gates in its constituent processing elements. (is that ironic or is it outrageous?)**

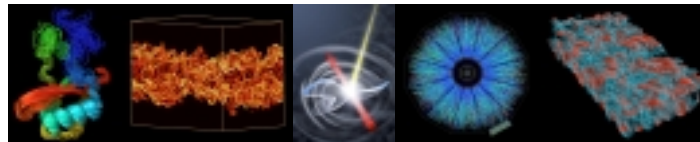


After John Shalf, NERSC



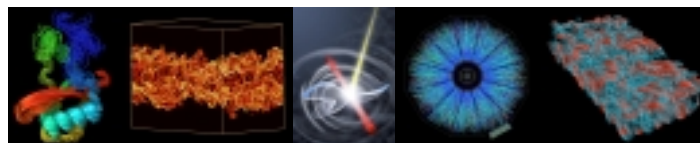
The complexity of a Petascale system is exceeding the complexity of its components

- Applications developers today write programs that are as complex as describing where every single bit must move between the 6,000 transistors of the 8080a.
- We need to at *least* get to the “assembly language” level.
- We may need to reconsider our entire hardware/software programming model if this is indeed what the future holds for us.



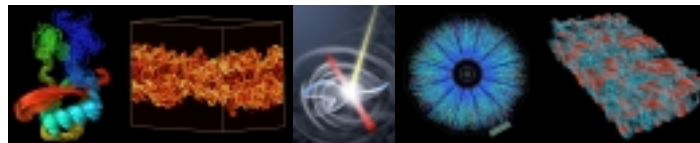
Overview

- History and Future of Petaflops Computing
- HPC in 2006: “It was the best of times, it was the worst of times ...”
- “A Petaflops before its Time”
- The power problem
- The scaling problem
- What's next?



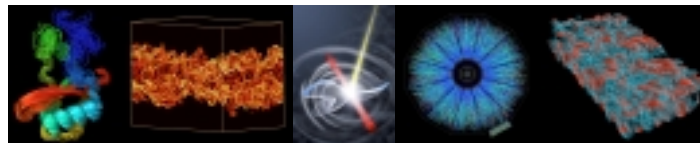
Until 2011: The best of times ...

- The high end technical market (departmental, server, workgroup) will continue to expand - clusters everywhere
- Scalability to a few 1000 processors easily feasible with current technologies
- Software opportunities
- Entry of big players like Microsoft



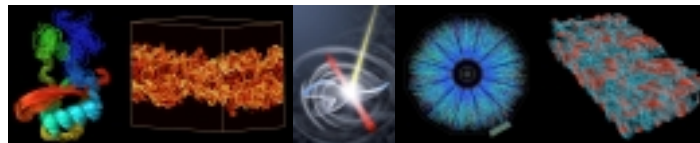
Until 2011: The worst of times ...

- The bad (cheap clusters) will drive out the good - further decline of the HPC capability market
- “A petaflop before its time” - early and easy successes will detract from solutions to the difficult problems of building Petascale systems



Ushering in True Petascale Computing: Challenge and Opportunity 2008 - 2016

- **All of computing** will be highly parallel by 2010
- The current ecosystem will become untenable after about 2010 in the face of the architectural, software, and power challenges
- Being on the forefront of these challenges, HPC has the opportunity to completely redefine computing
- How are we going to change the ecosystem?
- What are we going to change it into?



It was the best of times, it was the worst of times, it was the age of wisdom, it was the age of foolishness, it was the epoch of belief, it was the epoch of incredulity, it was the season of Light, it was the season of Darkness, it was the spring of hope, it was the winter of despair, we had everything before us, we had nothing before us, we were all going direct to heaven, we were all doing direct the other way - in short, the period was so far like the present period, that some of its noisiest authorities insisted on its being received, for good or for evil, in the superlative degree of comparison only.

Charles Dickens (1812 - 1870), A Tale of Two Cities

